

# Quantum Self-Supervised Learning

B. Jaderberg<sup>1</sup>, L. W. Anderson<sup>1</sup>, W. Xie<sup>2</sup>, S. Albanie<sup>2</sup>, M. Kiffner<sup>1,3</sup> and D. Jaksch<sup>1,3</sup>

<sup>1</sup>*Clarendon Laboratory, University of Oxford, Parks Road, Oxford OX1 3PU, United Kingdom*

<sup>2</sup>*Visual Geometry Group, Department of Engineering Science, University of Oxford and*

<sup>3</sup>*Centre for Quantum Technologies, National University of Singapore, 3 Science Drive 2, Singapore 117543*

(Dated: March 30, 2021)

The popularisation of neural networks has seen incredible advances in pattern recognition, driven by the supervised learning of human annotations. However, this approach is unsustainable in relation to the dramatically increasing size of real-world datasets. This has led to a resurgence in self-supervised learning, a paradigm whereby the model generates its own supervisory signal from the data. Here we propose a hybrid quantum-classical neural network architecture for contrastive self-supervised learning and test its effectiveness in proof-of-principle experiments. Interestingly, we observe a numerical advantage for the learning of visual representations using small-scale quantum neural networks over equivalently structured classical networks, even when the quantum circuits are sampled with only 100 shots. Furthermore, we apply our best quantum model to classify unseen images on the *ibmq-paris* quantum computer and find that current noisy devices can already achieve equal accuracy to the equivalent classical model on downstream tasks.

## I. INTRODUCTION

In the past decade, machine learning has revolutionised scientific analysis, yielding breakthrough results in protein folding [1], black hole imaging [2] and heart disease treatment [3]. At the forefront of this progress is deep learning [4], characterised by the successive application of artificial neural network layers [5, 6]. Notably, its use in computer vision has seen the top-1 accuracy on benchmark datasets such as ImageNet soar from 52% [7] to over 90% [8], fuelled by shifts in the underlying techniques used [9, 10]. However, what has remained consistent in these top performing models is the use of labelled data to supervise the representation learning process. Whilst effective, the reliance on large quantities of human-provided annotations presents a significant challenge as to whether such approaches will scale into the future. Crucially, modern datasets such as the billions of images uploaded to social media are both vast and unbounded in their subject, quickly making the task of labelling unfeasible.

This has reignited interest in an alternative approach, termed *self-supervised learning* [11], which seeks instead to exploit structure in the data itself as a learning signal. Rather than predict human annotations, a model is trained to perform a *proxy task*, that makes use of attributes of the data that can be inferred without labelling. Furthermore, the proxy task should encourage the model to learn representations that capture useful factors of variation in the visual input, such that solving it ultimately correlates with solving tasks of interest after training. Recent progress in the self-supervised learning of visual data has been driven by the success of contrastive learning [12–16], in which the proxy task is differentiating augmented instances of the same image from all other images. Provided the correct choice of augmentations, this produces a model which is invariant to transformations that do not change the semantic meaning of the image, allowing the learning of recognisable

features and patterns in unlabelled datasets.

With these techniques, contrastive learning is able to learn visual representations with comparable quality to supervised learning [16, 17], without the bottleneck of labelling. However, it is a fundamentally more difficult task than its supervised counterpart, and capturing complex correlations between augmented views requires more training data, more training time and larger network capacity [14, 15]. Therefore, it is important to consider whether emerging technologies can contribute to the growing requirement for more powerful neural networks. Variational quantum algorithms (VQAs) [18], a near term application of quantum computing, are one such new paradigm. While VQAs have been used to solve many types of optimisation problems [19–23], it is their application to supervised learning [24–26], unsupervised learning [27], generative models [28, 29] and reinforcement learning [30–32] which has led to them being referred to as quantum neural networks (QNNs) [33–35]. Importantly, although it has not yet been shown generally [36], early evidence suggests that QNNs can achieve an advantage over their classical counterparts for specific problems [37], driven by access to an exponentially larger quantum feature space [25]. Furthermore, the findings that QNNs are able to achieve a higher effective dimension than classical neural networks on real-world data [38] provides a theoretical framework in which quantum computers can provide an advantage for deep learning tasks. In this work, we construct a contrastive learning architecture in which classical and quantum neural networks are trained together. By randomly augmenting each image in the dataset, our hybrid network learns visual representations which groups different views of the same image together in both classical and Hilbert space. Afterwards, we test the quality of the representations by using them to train a linear classifier, which then makes predictions on an unseen test set. We find that our hybrid encoder, constrained in both size and training time by quantum simulation overheads, achieves an average test accuracy

of  $(46.51 \pm 1.37)\%$ . In contrast, replacing the QNN with a classical neural network of equivalent width and depth results in a model which obtains  $(43.49 \pm 1.31)\%$  accuracy. Thus, our results provide the first indication that a quantum model may better capture the complex correlations required for self-supervised learning.

We then apply the best performing quantum model to classify test images on a real quantum computer. Notably, the accuracy achieved using the *ibmq-paris* [39] device equals the best performing classical model, despite significant device noise. This illustrates the capability of our algorithm for real-world applications using current devices, with flexibility to assign more of the encoding to QNNs as quantum hardware improves. While further research is required to demonstrate scalability, our scheme provides a strong foundation for quantum self-supervised learning. Excitingly, given that contrastive learning has also been successfully applied to non-visual data [15, 40–43], our work opens the possibility of using QNNs to learn large, unlabelled datasets across a range of disciplines.

## II. RESULTS

### A. Contrastive learning architecture

Given an unlabelled dataset, the objective of self-supervised learning is to find low dimensional encodings of the images which retain important higher level features. In this work, we train a model to do this by adapting the widely used SimCLR algorithm [17], the steps of which can be seen in Fig. 1. Firstly for a given image, the data of which is contained within  $\tilde{x}_i$ , we generate two augmentation functions. Each one randomly crops, rotates, blurs and colour distorts the picture, such that two augmented views  $\tilde{x}_i^1, \tilde{x}_i^2$  of the same base image are produced. Importantly, these augmentations still allow for the underlying object to remain visually distinguishable. This enables us to assert that these two views contain a recognisable description of the same class, which we call a positive pair.

Once this positive pair is generated, each view is passed through a set of neural networks. First, an encoder network is applied, which maps the high dimensional input data  $\tilde{x}_i^1, \tilde{x}_i^2$  to low dimensional representations  $\tilde{y}_i^1, \tilde{y}_i^2$ . Then the output of the encoder network is passed to the projection head, a small multi-layer-perceptron (MLP) [44] consisting of two fully connected layers. This produces the final representations  $\tilde{z}_i^1, \tilde{z}_i^2$ .

Given a batch of  $N$  images, the above process is repeated such that we are left with  $2N$  representations corresponding to  $2N$  augmented views. Looking at all possible pairings of these representations, we have not only positive pairs (e.g.,  $\tilde{z}_i^1, \tilde{z}_i^2$ ) but also negative pairs (e.g.,  $\tilde{z}_i^1, \tilde{z}_j^1$  where  $i \neq j$ ), which we cannot definitely say contain the same class. For each training step, all of these possible pairs are used to calculate the normalised temperature-scaled cross entropy loss (NT-Xent) [45] (see Methods),

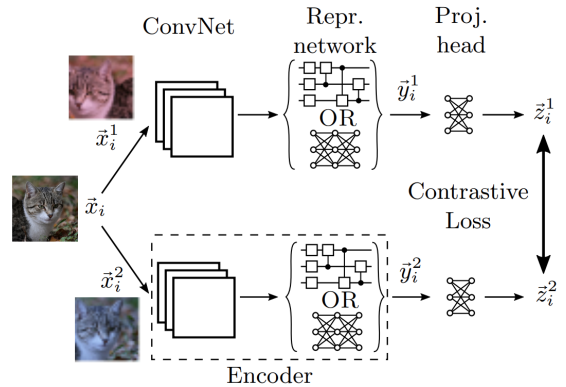


FIG. 1. Schematic of the overall neural network architecture and contrastive training method. For each input image  $\tilde{x}_i$ , a pair of random augmentations are generated and applied to form a positive pair  $\tilde{x}_i^1, \tilde{x}_i^2$ . These are transformed by the encoder network, consisting of classical convolutional layers and a quantum or classical representation network, into representation vectors  $\tilde{y}_i^1, \tilde{y}_i^2$ . The projection head subsequently maps the representations to the vectors  $\tilde{z}_i^1, \tilde{z}_i^2$ , such that contrastive loss can be applied without inducing loss of information on the encoder.

which is minimised via stochastic gradient descent [46]. Intuitively, minimising this loss function can be understood as training the network to produce representations in which positive pairs are mapped close together and negative pairs far apart, as measured by their cosine similarity. This idea is a core concept in contrastive learning and many machine learning techniques [47]. Note that whilst it is possible to train the network by applying NT-Xent directly to the output of the encoder, the contrastive loss function is known to induce loss of information on the layer it is applied to [17]. Therefore, the addition of the projection head ensures that the encoder remains sensitive to image characteristics (e.g., colour, orientation) that improves performance on downstream tasks.

In order to incorporate QNNs, we modify the encoder to contain both classical and quantum layers working together. The first part of the encoder consists of a convolutional neural network, which in this work is the widely used ResNet-18 [48]. This produces a 512 length feature vector, which is already an initial encoding of the augmented image. However, we then extend the encoder with a second network, which we call the representation network as it acts directly on the representation space. This consists of either a multi-layer QNN of width  $W$ , or a classical fully connected MLP with equivalent width and depth. Ideally the representation network would have width  $W = 512$ , so as to minimise loss of information. However, we instead look to work in a regime which is realisable on current quantum computers, and as such in this work we use  $W = 8$ . This is achieved by following the convolutional network with a single classical layer that compresses the vector, a common technique used to link classical and quantum networks together [49, 50].

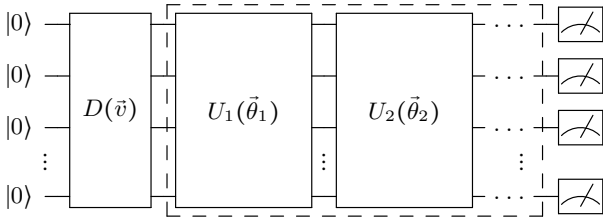


FIG. 2. General structure of a QNN. An input vector  $\vec{v}$  is encoded into the qubits by a data loading unitary  $\hat{D}(\vec{v})$ . The variational ansatz consists of layers  $\{\hat{U}_1(\vec{\theta}_1), \hat{U}_2(\vec{\theta}_2), \dots\}$  and is parameterised by trainable parameters  $\{\vec{\theta}_1, \vec{\theta}_2, \dots\}$ . The output of the QNN is taken as the average of repeated measurements in the  $\hat{\sigma}_z$  basis.

After the representation network is applied, the resultant encoding is passed onto the previously described projection head. To maintain the structure of the original SimCLR architecture, we limit the projection head to be no wider than the width of the QNN.

### B. Quantum representation network

The quantum representation network follows the structure shown in Fig. 2, beginning with a data loading unitary  $\hat{D}(\vec{v})$ . Whilst schemes exist to encode data into quantum circuits with exponential compression [51, 52], these require a prohibitively large number of logic gates compared to current hardware capabilities. By compressing the output of the ConvNet as described in the previous section, we need only to solve the simpler issue of loading a vector  $\vec{v}$  of length  $W$  into equally as many qubits. This is achieved by applying a single qubit rotation  $\hat{R}_x$  to each qubit in the register;  $\hat{D}(\vec{v}) = \bigotimes_{i=1}^W \hat{R}_x(v^k)$ . Here,  $v^k$  is the  $k$ th element of input vector  $\vec{v}$  and is mapped to the range  $[0, \pi]$  to prevent large values wrapping back around the Bloch sphere.

Once the input data is loaded, we apply the learning component of our QNN, a parameterised quantum circuit ansatz. In applications where the ansatz is used to solve optimisation problems relating to a physical system (e.g., the simulation of molecules), the circuit structure and choice of logic gates can be inspired by the underlying Hamiltonian [53]. However, without such symmetries to guide our choice, we use a variational ansatz based on recent theoretical findings in expressibility and entangling capability [54]. The ansatz is shown in Fig. 3, the structure of which is derived from circuit 14 of Ref. [54] and was chosen due to its performance in both these metrics. After the application of several ansatz layers, the network is finished by measuring each qubit to obtain an expectation value in the  $\hat{\sigma}_z$  basis. When evaluated on a real quantum computer or sampling-based simulator, the expectation value is constructed by averaging the sampled eigenvalues over a finite number of shots. If evaluated on

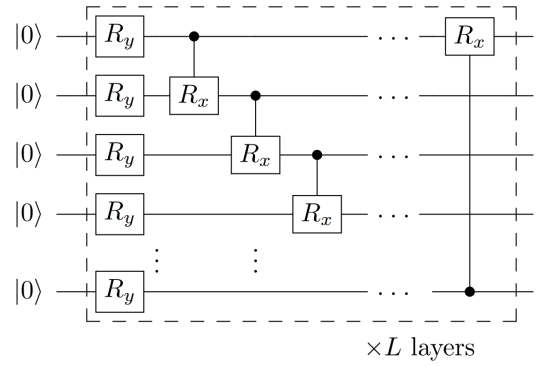


FIG. 3. (a) Variational ansatz used in this work. Each layer consists of a single qubit  $\hat{R}_y$  rotation on each qubit, followed by controlled  $\hat{R}_x$  rotations, connecting the qubits in a ring topology. Every rotation gate is parameterised by a different variational parameter.

a statevector simulator, the expectation value is calculated exactly.

The gradients of the QNN output with respect to the trainable parameters and the input parameters are calculated using the parameter shift rule [33, 55], which we describe here. Consider an observable  $\hat{O}$  measured on the state

$$|\psi(\vec{\theta})\rangle = \prod_i \hat{U}_i(\theta_i) \hat{V}_i |0\rangle, \quad (1)$$

resulting from the application of  $M$  parameterised gates  $\hat{U}_1, \hat{U}_2, \dots, \hat{U}_M$  and  $M$  fixed gates  $\hat{V}_1, \hat{V}_2, \dots, \hat{V}_M$ , where gates  $\hat{U}_i = e^{i\theta_i \hat{P}_i/2}$  are generated by operators  $\hat{P}_i \in \{1, \hat{\sigma}_x, \hat{\sigma}_y, \hat{\sigma}_z\}^{\otimes n}$  that are tensor products of the Pauli operators. According to the parameter shift rule, the gradient of the expectation value  $f = \langle \psi(\vec{\theta}) | \hat{O} | \psi(\vec{\theta}) \rangle$  with respect to parameter  $\theta_i$  is given by

$$\frac{\partial f(\vec{\theta})}{\partial \theta_i} = \frac{1}{2} \left[ f\left(\theta_i + \frac{\pi}{2}\right) - f\left(\theta_i - \frac{\pi}{2}\right) \right]. \quad (2)$$

For each parameterised gate within the circuit, including both the variational ansatz and data loading unitary, an unbiased estimator for the gradient is calculated by measuring the QNN with the two shifted parameter values given in Eq. (2).

Once the QNN gradients have been calculated, we combine them with gradients of the classical components to obtain gradients of the loss function with respect to all trainable quantum and classical parameters via back-propagation [56]. In this way, the QNN is trained simultaneously with the classical networks, and the quality of the gradients produced on quantum hardware play a crucial role in the training ability of the whole network.

### C. Training

To examine whether the proposed architecture can successfully train, we apply it to the CIFAR-10 dataset [57].

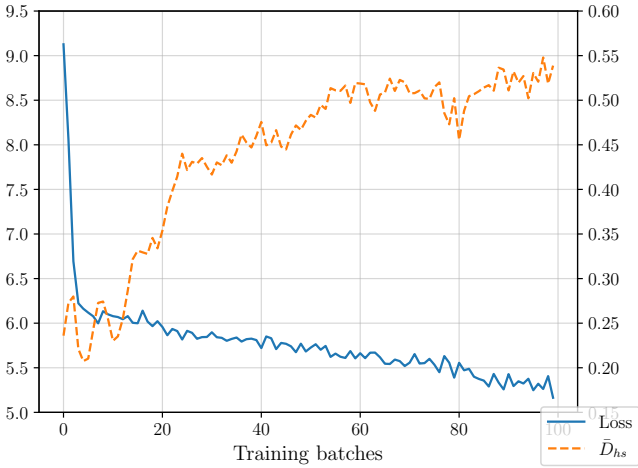


FIG. 4. Contrastive learning with a quantum representation network. After each batch of 256 images, the loss function (blue, left axis) and average Hilbert Schmidt distance between positive pairs and negative pairs (orange, right axis) is recorded.

In this preliminary experiment we restrict the dataset to the first two classes, leaving 10,000  $32 \times 32$  colour images containing either an aeroplane or automobile. We also train this initial model without a projection head, since it is not being used for classification later. The classical and quantum components are trained together from scratch, achieved by integrating the Qiskit [58] and PyTorch [59] frameworks together. The full list of training hyperparameters can be found in the Methods section.

Fig. 4 shows the results of two key metrics after training for 100 batches. Firstly, we record the loss after each batch, the minimisation of which represents the ability to produce representations in the classical  $W$  dimensional space whereby positive pairs have high similarity. Our results show that the loss decreases from 9.13 to 5.16 over the course of training, indicating that our model is able to learn. Importantly, since the quantum and classical parameters are trained together, this shows that information is successfully passed both forwards and backwards between these different network paradigms.

Secondly, we log the Hilbert-Schmidt distance ( $D_{\text{HS}}$ ), a metric that has been applied in quantum machine learning previously to study data embedding in Hilbert space [50]. Here, we use it to track the separation between our pseudo classes in the  $2^W$  dimensional quantum state space while optimising the classical loss function. For a given positive pair  $\tilde{x}_i^1, \tilde{x}_i^2$ , we calculate the statistical ensembles

$$\rho_i = \frac{1}{2} (|\psi_i^1\rangle\langle\psi_i^1| + |\psi_i^2\rangle\langle\psi_i^2|), \quad (3a)$$

$$\sigma_i = \frac{1}{2N-2} \sum_{j \neq i} (|\psi_j^1\rangle\langle\psi_j^1| + |\psi_j^2\rangle\langle\psi_j^2|), \quad (3b)$$

where  $|\psi_i^\alpha\rangle$  is the statevector produced by the hybrid en-

coder given augmented view  $\tilde{x}_i^\alpha$ . The Hilbert-Schmidt distance is then given by

$$D_{\text{HS},i} = \text{tr}((\rho_i - \sigma_i)^2). \quad (4)$$

We repeat this for each positive pair in the batch and record the mean,  $\bar{D}_{\text{HS}} = \frac{1}{N} \sum_i D_{\text{HS},i}$ . From Fig. 4, we see that  $\bar{D}_{\text{HS}}$  increases consistently across the range of training. This indicates that the QNN successfully learns to group positive pairs and separate them from other images in Hilbert space. Thus, we show that the quantum component of the encoder contributes to the overall learning process, despite the network's parameters being optimised explicitly in a classical space.

#### D. Linear probing

Once training is complete, we require a way to test the quality of the image representations learnt by the encoder. Specifically, a good encoding will produce representations whereby different classes are linearly separable in the representation space [60]. Therefore, we numerically test the encoder using the established linear evaluation protocol [60], in which a linear classifier is trained on the output of the encoder network, whilst the encoder is frozen to stop it training any further. Once this linear probe experiment has trained for 100 epochs, we apply the whole network to unseen test data and record the classification accuracy.

#### E. Quantum and classical results on the simulator

We repeat training, this time with the first five classes of CIFAR-10 and a projection head. We train models with three different types of representation networks; classical MLP with bias and Leaky ReLU activation functions after each layer, quantum trained on a statevector simulator and quantum trained on a sampling-based simulator. We choose the representation networks to be width  $W = 8$  in order to minimise the simulation overhead, whilst still being in a compression regime where training is stable (see Appendix A). Quantitatively, this means our two-layer classical and quantum representation networks have 144 and 32 learnable parameters respectively.

Fig. 5 shows the result of linear probe experiments at checkpoints across 176 batches of contrastive training. We find that when the quantum circuits are evaluated using a statevector simulator, the quantum representation network produces higher average accuracy on the test set than the equivalent classical network at all points probed throughout training, and is separated by more than one standard deviation for over half of these. The highest accuracy is obtained at the end of training, where the quantum model achieves an accuracy of  $(46.51 \pm 1.37)\%$  compared to  $(43.49 \pm 1.31)\%$  for the classical model. We

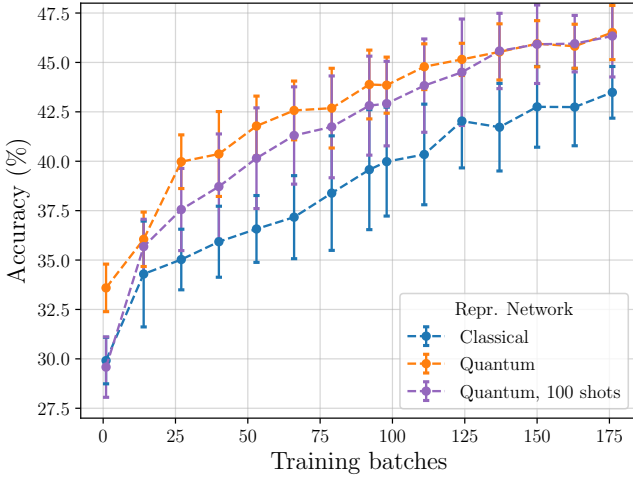


FIG. 5. Classification accuracy achieved in linear probing experiments using the encoder at checkpoints across self-supervised training. Comparison between models trained with a classical representation network (blue), quantum representation network evaluated on a statevector simulator (orange) and quantum representation network evaluated on a sampling-based simulator with 100 shots (purple). The markers show the average of six independently trained models, whilst the error bars show one standard deviation.

also find that this numerical advantage is dependent on the correct choice of ansatz (see Appendix B).

Subsequently, we explore whether using a finite number of shots limits this advantage. We train another quantum model on a simulator where the expectation values of measured qubits are sampled from 100 shots, both in the forward pass (generating the representations) as well as the backwards pass (calculating gradients). We find that beyond the first batch, the average accuracy of this model is still above what is achieved by the classical representation network, reaching  $(46.34 \pm 2.07)\%$  by the end of training. Significantly, this matches the performance of the statevector simulator, which represents the limit of infinite shots, demonstrating resilience of our scheme to shot noise. However, we note that the additional uncertainty introduced by the sampling does manifest as a larger standard deviation between repeated runs, compromising the consistency of the advantage.

## F. Real device experiments

In the previous section, we showed that a numerical advantage can be achieved for self-supervised learning with a quantum representation network, even when sampling the quantum circuits with only 100 shots. However, it does not follow that such an improvement can necessarily be realised on current quantum devices. The biggest barrier to this is the complex noise present on quantum hardware, a product of both the finite lifetime that qubits can be held in coherent states for and imperfections in

		Predicted									
		aero.	auto.	bird	cat	deer	aero.	auto.	bird	cat	deer
True	aeroplane	123	27	15	3	5	124	26	8	12	3
	automobile	34	111	13	21	3	43	101	6	25	7
	bird	17	21	28	40	70	18	13	36	53	56
	cat	2	25	29	75	50	6	32	15	87	41
	deer	11	21	23	54	79	16	19	28	50	75

(a)

(b)

(a)

(b)

FIG. 6. Confusion matrix from classifying 900 images using the best performing (a) classical model evaluated on a classical computer (b) quantum model evaluated on a real quantum computer with 100 shots per circuit. For a given true label (rows) and predicted label (columns), the number in each box shows the total number of times that prediction was made.

the application of logic gates. To this end, we test the ability of real devices to accurately prepare representations produced by a pretrained quantum model and how this changes downstream accuracy on the test set.

We construct a linear probe experiment with a quantum representation network and load in weights from the best performing pretrained model in which circuits were evaluated with 100 shots. Freezing all of the layers so that the entire network no longer trains, we repeat classification of images from the test set, however this time the circuits are executed on IBM's 27-qubit *ibmq-paris* quantum computer. To reduce the number of gates, particularly SWAP operations caused by a mismatch between the ansatz and physical qubit connectivities, the circuits are recompiled using incremental structural learning [61] before execution (see Methods).

Fig. 6a shows the result of classifying 900 images randomly sampled from the test set, using the best performing classical model and evaluated on a classical computer. Fig. 6b shows the result when classifying the same images using the best performing 100-shot quantum model, evaluated on *ibmq-paris*. Overall, the classical and quantum models achieve an accuracy of 47.27% and 47.00% respectively. Excitingly, this demonstrates that in this experiment, error induced by noise on the quantum computer is able to be offset by the enhanced theoretical performance of quantum neural networks, provided the circuit depth is reduced with recompilation techniques. Furthermore, in both setups the most correctly predicted class was aeroplanes (71.1% and 71.7%) whilst the most incorrectly predicted class was birds (15.9% and 20.5%), both of which the quantum model performed better on. We propose that birds and deer were most likely to be mistaken with one another due to the images sharing a common background of the outdoor natural environment.



### III. DISCUSSION

In this work, we propose a hybrid quantum-classical architecture for self-supervised learning and demonstrate a numerical advantage in the learning of visual representations using small-scale QNNs. We train quantum and classical neural networks together, such that encodings are learnt that maximise the similarity of augmented views of the same image in the representation space, as well as implicitly in Hilbert space. After training is complete, we determine the quality of the embedding by tasking a linear probe to classify images from different classes. We find that an encoder with a QNN acting in the representation space achieves higher average test set accuracy than a classical neural network with equivalent width and depth, even when evaluating quantum circuits with only 100 shots.

We then apply our best performing pretrained classical and quantum models to downstream classification, whereby the quantum circuits were evaluated on a real quantum computer. The observation of a quantum predictive signal with equivalent accuracy to that of the classical model, despite the complex noise present on current quantum devices, is representative of the potential practical benefit of our setup. If recent progress in superconducting qubit hardware continues [62–64], it is likely that QNNs running on real devices will outperform equally sized classical neural networks in the near future in this experiment.

One advantage of the hybrid approach taken in this paper is the resulting flexibility in how much of the encoder is quantum or classical. As the quality and size of quantum hardware improves, our scheme allows classical capacity to be substituted for quantum, eventually replacing ResNet entirely. By optimising directly for the Hilbert-Schmidt distance, it is possible with a fully quantum encoder to apply our setup to problems in which the data is itself quantum [65–67]. Promisingly, in this regime it may prove that the advantage observed in this work is further extended, given the ability of a quantum model to inherently exploit the dimensionality of the input [68]. With classical contrastive learning having been applied to non-visual problems in biology [42] and chemistry [43], our work provides a strong foundation for applying quantum self-supervised learning to fundamentally quantum problems in the natural sciences [69].

Looking forward, an open question remains as to whether a general quantum advantage for self-supervised learning may prove possible [62, 70], in which no classical computer of any size can produce accuracies equal to that of a quantum model. Achieving this would likely require a QNN with width greater than 60 qubits, such that the dimensionality of the accessible feature space becomes classically intractable. Therefore, considerable research still remains into the scalability of our scheme, which was only demonstrated at scales feasible on current quantum hardware. Furthermore, training the QNNs on real devices remains a challenge, due to both the lack of dedi-

cated access to execute millions of circuits and significant gate errors. Nevertheless, further research is warranted into the potential noise tolerance of our algorithm due to its unique ability to use random transformations as a learning signal. This means that it may be possible to obtain accurate results in the presence of noise, so long as the errors act perpendicular to the semantic meaning of the data. In this way, quantum self-supervised learning and the scheme put forward in this work may yet prove a candidate for general quantum advantage on noisy devices.

### IV. METHODS

#### A. Contrastive Loss Function

Here we formally define the process of contrastive learning. Let us have an augmentation function  $\xi(\cdot; a)$ . This augmentation combines cropping, rotation, Gaussian blurring and colour distortion of the image, and the amount by which each of these operations is performed is governed by a list of continuous random variables  $a$ . Each time we apply an augmentation, we randomly sample  $a$  from a distribution  $A$  such that applications of the augmentation function are independent from one another. Given a batch  $\mathcal{B} = \{\bar{x}_1, \bar{x}_2, \dots, \bar{x}_N\}$  we apply  $\xi(\cdot; a)$  twice to each image. For a particular image  $\bar{x}_i$ , we now have a pair of views  $\mathcal{P}_i = \{\xi(\bar{x}_i; a_1), \xi(\bar{x}_i; a_2) \mid a_1, a_2 \sim A\}$  which came from the same base image. We call this a positive pair. Conversely, we define the negative set  $\mathcal{N}_i = \{\xi(\bar{x}_n; a_{n,1}), \xi(\bar{x}_n; a_{n,2}) \mid a_{n,1}, a_{n,2} \sim A; \forall n \neq i\}$  of  $2N - 2$  elements which contains a randomly sampled pair of augmentations for each of the  $N - 1$  remaining images. During contrastive training, all augmented views within the batch are passed through our architecture. The encoder network  $f(\cdot) : \bar{x} \rightarrow \vec{y}$  and projection head  $g(\cdot) : \vec{y} \rightarrow \vec{z}$ , are applied to give outputs  $\vec{z}_i^\alpha = g(f(\xi(\bar{x}_i; a)))$ ;  $a \sim A$  for each of the two arms (labelled by  $\alpha = 1, 2$ ). Using this batch of representations, the NT-Xent loss is calculated as:

$$\mathcal{L} = \sum_{\substack{i \in \{1, \dots, N\} \\ \{\bar{x}_i^1, \bar{x}_i^2\} = \mathcal{P}_i}} \log \frac{-\exp(\vec{z}_i^1 \cdot \vec{z}_i^2 / \tau)}{\exp(\vec{z}_i^1 \cdot \vec{z}_i^2 / \tau) + \sum_{\bar{x}_n^\alpha \in \mathcal{N}_i} \exp(\vec{z}_i^1 \cdot \vec{z}_n^\alpha / \tau)}, \quad (5)$$

where  $\vec{z}_i^1 \cdot \vec{z}_i^2$  is the cosine similarity between representation vectors corresponding to the first and second random augmentations applied to  $\bar{x}_i$ . The product  $\vec{z}_i^1 \cdot \vec{z}_{n \neq i}^\alpha$  is defined similarly for a negative pair. The summation covers all augmented inputs  $\bar{x}_i^1$  and the corresponding positive pair  $\bar{x}_i^2$ . The summation is over both orderings  $\{\bar{x}_i^1, \bar{x}_i^2\} = \mathcal{P}_i$  and  $\{\bar{x}_i^2, \bar{x}_i^1\} = \mathcal{P}_i$  giving a total of  $2N$  terms. The softmax temperature  $\tau$  is a hyperparameter used to control how conservative the model is. Overall, optimising this loss can be seen as a task of instance discrimination, in which we produce a network which emits

vectors with higher similarity scores between the augmented views of the same instance (positive pairs) than augmented views of other instances (negative pairs).

## B. Training hyperparameters

Throughout this work, the training parameters used are; batch size: 256, optimiser: ADAM [71] with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ , learning rate:  $10^{-3}$ , weight decay:  $10^{-6}$  and softmax temperature: 0.07.

## C. Recompilation of quantum neural networks

When executing QNNs on the *ibmq-paris* device, translating the ring topology of our variational ansatz to the honeycomb structure that the qubits are physically connected by requires a significant number of SWAP operations. Quantitatively this increases the number of two-qubit gates in the circuit from 16 to 143, which poses a significant challenge to obtaining a predictive signal beyond random noise since the total circuit error scales exponentially with the number of gates. To mitigate this, for each image evaluated we approximately recompile the QNN using incremental structural learning (ISL) [61],

adapted so that only two-qubit connections available on the real device can be applied. Using this method, for over half of the executed circuits, an equivalent circuit is found which produces the same statevector with at least 99% overlap using on average 14 CNOT gates. For the remaining images, we apply ISL once again, but this time without any constraints on the connectivity of the circuit. This produces a shallower equivalent circuit with at least 99% overlap using on average 8 CNOT gates. Although some of these two qubit gates require SWAPs when implemented on the real device, they still represent a significant reduction in the depth of the circuit and total error incurred.

## ACKNOWLEDGMENTS

BJ, LWA, MK, DJ acknowledge support from the EPSRC National Quantum Technology Hub in Networked Quantum Information Technology (EP/M013243/1) and the EPSRC Hub in Quantum Computing and Simulation (EP/T001062/1). MK and DJ acknowledge financial support from the National Research Foundation, Prime Ministers Office, Singapore, and the Ministry of Education, Singapore, under the Research Centres of Excellence program. WX and SA are supported by EPSRC grant Seebibyte (EP/M013774/1) and Visual AI (EP/T028572/1).

- 
- [1] Andrew W Senior, Richard Evans, John Jumper, James Kirkpatrick, Laurent Sifre, Tim Green, Chongli Qin, Augustin Židek, Alexander WR Nelson, Alex Bridgland, *et al.*, “Improved protein structure prediction using potentials from deep learning,” *Nature* **577**, 706–710 (2020).
  - [2] K. Akiyama *et al.*, “First m87 event horizon telescope results. i. the shadow of the supermassive black hole,” *The Astrophysical Journal* **875**, 1–17 (2019).
  - [3] Christina V Theodoris, Ping Zhou, Lei Liu, Yu Zhang, Tomohiro Nishino, Yu Huang, Aleksandra Kostina, Sanjeev S Ranade, Casey A Gifford, Vladimir Uspenskiy, *et al.*, “Network-based screen in ipsc-derived cells reveals therapeutic candidate for heart valve disease,” *Science* **371** (2021).
  - [4] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton, “Deep learning,” *Nature* **521**, 436–444 (2015).
  - [5] Warren S McCulloch and Walter Pitts, “A logical calculus of the ideas immanent in nervous activity,” *The bulletin of mathematical biophysics* **5**, 115–133 (1943).
  - [6] Anders Krogh, “What are artificial neural networks?” *Nature biotechnology* **26**, 195–197 (2008).
  - [7] Yuanqing Lin, Fengjun Lv, Shenghuo Zhu, Ming Yang, Timothee Cour, Kai Yu, Liangliang Cao, and Thomas Huang, “Large-scale image classification: fast feature extraction and svm training,” in *CVPR 2011* (IEEE, 2011) pp. 1689–1696.
  - [8] Hieu Pham, Qizhe Xie, Zihang Dai, and Quoc V Le, “Meta pseudo labels,” *arXiv preprint arXiv:2003.10580* (2020).
  - [9] David G Lowe, “Object recognition from local scale-invariant features,” in *Proceedings of the seventh IEEE international conference on computer vision*, Vol. 2 (Ieee, 1999) pp. 1150–1157.
  - [10] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton, “Imagenet classification with deep convolutional neural networks,” *Advances in neural information processing systems* **25**, 1097–1105 (2012).
  - [11] Virginia R de Sa, “Learning classification with unlabeled data,” in *Advances in neural information processing systems* (1994) pp. 112–119.
  - [12] Zhirong Wu, Alexei A Efros, and Stella X Yu, “Improving generalization via scalable neighborhood component analysis,” in *Proceedings of the European Conference on Computer Vision (ECCV)* (2018) pp. 685–701.
  - [13] Zhirong Wu, Yuanjun Xiong, Stella Yu, and Dahua Lin, “Unsupervised feature learning via non-parametric instance-level discrimination,” *arXiv preprint arXiv:1805.01978* (2018).
  - [14] Olivier Henaff, “Data-efficient image recognition with contrastive predictive coding,” in *International Conference on Machine Learning* (PMLR, 2020) pp. 4182–4192.
  - [15] Aaron van den Oord, Yazhe Li, and Oriol Vinyals, “Representation learning with contrastive predictive coding,” *arXiv preprint arXiv:1807.03748* (2018).
  - [16] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick, “Momentum contrast for unsupervised visual representation learning,” in *Proceedings of the*

- IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2020) pp. 9729–9738.
- [17] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton, “A simple framework for contrastive learning of visual representations,” (2020), arXiv:2002.05709 [cs.LG].
  - [18] Marco Cerezo, Andrew Arrasmith, Ryan Babbush, Simon C Benjamin, Suguru Endo, Keisuke Fujii, Jarrod R McClean, Kosuke Mitarai, Xiao Yuan, Lukasz Cincio, *et al.*, “Variational quantum algorithms,” arXiv preprint arXiv:2012.09265 (2020).
  - [19] Alberto Peruzzo, Jarrod McClean, Peter Shadbolt, Man-Hong Yung, Xiao-Qi Zhou, Peter J Love, Alán Aspuru-Guzik, and Jeremy L O’Brien, “A variational eigenvalue solver on a photonic quantum processor,” *Nature communications* **5**, 4213 (2014).
  - [20] Google AI Quantum *et al.*, “Hartree-fock on a superconducting qubit quantum computer,” *Science* **369**, 1084–1089 (2020).
  - [21] Edward Farhi, Jeffrey Goldstone, and Sam Gutmann, “A quantum approximate optimization algorithm,” arXiv preprint arXiv:1411.4028 (2014).
  - [22] Leo Zhou, Sheng-Tao Wang, Soonwon Choi, Hannes Pichler, and Mikhail D Lukin, “Quantum approximate optimization algorithm: performance, mechanism, and implementation on near-term devices,” arXiv preprint arXiv:1812.01041 (2018).
  - [23] He Ma, Marco Govoni, and Giulia Galli, “Quantum simulations of materials on near-term quantum computers,” *npj Computational Materials* **6**, 1–8 (2020).
  - [24] Edward Grant, Marcello Benedetti, Shuxiang Cao, Andrew Hallam, Joshua Lockhart, Vid Stojevic, Andrew G. Green, and Simone Severini, “Hierarchical quantum classifiers,” *npj Quantum Information* **4**, 65 (2018).
  - [25] Vojtěch Havlíček, Antonio D. Córcoles, Kristan Temme, Aram W. Harrow, Abhinav Kandala, Jerry M. Chow, and Jay M. Gambetta, “Supervised learning with quantum-enhanced feature spaces,” *Nature* **567**, 209–212 (2019).
  - [26] Maria Schuld, Alex Bocharov, Krysta M. Svore, and Nathan Wiebe, “Circuit-centric quantum classifiers,” *Phys. Rev. A* **101**, 032308 (2020).
  - [27] JS Otterbach, R Manenti, N Alidoust, A Bestwick, M Block, B Bloom, S Caldwell, N Didier, E Schuyler Fried, S Hong, *et al.*, “Unsupervised machine learning on a hybrid quantum computer,” arXiv preprint arXiv:1712.05771 (2017).
  - [28] Marcello Benedetti, Delfina Garcia-Pintos, Oscar Perdomo, Vicente Leyton-Ortega, Yunseong Nam, and Alejandro Perdomo-Ortiz, “A generative modeling approach for benchmarking and training shallow quantum circuits,” *npj Quantum Information* **5**, 1–9 (2019).
  - [29] Christa Zoufal, Aurélien Lucchi, and Stefan Woerner, “Quantum generative adversarial networks for learning and loading random distributions,” *npj Quantum Information* **5**, 1–9 (2019).
  - [30] Samuel Yen-Chi Chen, Chao-Han Huck Yang, Jun Qi, Pin-Yu Chen, Xiaoli Ma, and Hsi-Sheng Goan, “Variational quantum circuits for deep reinforcement learning,” *IEEE Access* **8**, 141007–141024 (2020).
  - [31] Owen Lockwood and Mei Si, “Reinforcement learning with quantum variational circuit,” in *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, Vol. 16 (2020) pp. 245–251.
  - [32] V. Saggio, B. E. Asenbeck, A. Hamann, T. Strömberg, P. Schiansky, V. Dunjko, N. Friis, N. C. Harris, M. Hochberg, D. Englund, and *et al.*, “Experimental quantum speed-up in reinforcement learning agents,” *Nature* **591**, 229–233 (2021).
  - [33] K. Mitarai, M. Negoro, M. Kitagawa, and K. Fujii, “Quantum circuit learning,” *Physical Review A* **98** (2018), 10.1103/physreva.98.032309.
  - [34] Kerstin Beer, Dmytro Bondarenko, Terry Farrelly, Tobias J. Osborne, Robert Salzmann, Daniel Scheiermann, and Ramona Wolf, “Training deep quantum neural networks,” *Nature Communications* **11**, 808 (2020).
  - [35] Marcello Benedetti, Erika Lloyd, Stefan Sack, and Mattia Fiorentini, “Parameterized quantum circuits as machine learning models,” *Quantum Science and Technology* **4**, 043001 (2019).
  - [36] Logan G. Wright and Peter L. McMahon, “The capacity of quantum neural networks,” (2019), arXiv:1908.01364 [quant-ph].
  - [37] Hsin-Yuan Huang, Michael Broughton, Masoud Mohseni, Ryan Babbush, Sergio Boixo, Hartmut Neven, and Jarrod R McClean, “Power of data in quantum machine learning,” arXiv preprint arXiv:2011.01938 (2020).
  - [38] Amira Abbas, David Sutter, Christa Zoufal, Aurélien Lucchi, Alessio Figalli, and Stefan Woerner, “The power of quantum neural networks,” arXiv preprint arXiv:2011.00027 (2020).
  - [39] IBM Quantum. <https://quantum-computing.ibm.com/>, 2021.
  - [40] Andriy Mnih and Koray Kavukcuoglu, “Learning word embeddings efficiently with noise-contrastive estimation,” in *Conference on Neural Information Processing Systems* (2013).
  - [41] Aditya Grover and Jure Leskovec, “Node2vec: Scalable feature learning for networks,” in *ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD)* (2016).
  - [42] Amy X Lu, Haoran Zhang, Marzyeh Ghassemi, and Alan M Moses, “Self-supervised contrastive learning of protein representations by mutual information maximization,” *bioRxiv* (2020).
  - [43] Sabrina Jaeger, Simone Fulle, and Samo Turk, “Mol2vec: unsupervised machine learning approach with chemical intuition,” *Journal of chemical information and modeling* **58**, 27–35 (2018).
  - [44] Ke-Lin Du and Madisetti NS Swamy, *Neural networks and statistical learning* (Springer Science & Business Media, 2013).
  - [45] Kihyuk Sohn, “Improved deep metric learning with multi-class n-pair loss objective,” in *Proceedings of the 30th International Conference on Neural Information Processing Systems* (2016) pp. 1857–1865.
  - [46] Herbert Robbins and Sutton Monroe, “A stochastic approximation method,” *The annals of mathematical statistics*, 400–407 (1951).
  - [47] Raia Hadsell, Sumit Chopra, and Yann LeCun, “Dimensionality reduction by learning an invariant mapping,” in *2006 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’06)*, Vol. 2 (IEEE, 2006) pp. 1735–1742.
  - [48] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition* (2016) pp. 770–778.



- [49] Andrea Mari, Thomas R. Bromley, Josh Izaac, Maria Schuld, and Nathan Killoran, “Transfer learning in hybrid classical-quantum neural networks,” *Quantum* **4**, 340 (2020).
- [50] Seth Lloyd, Maria Schuld, Aroosa Ijaz, Josh Izaac, and Nathan Killoran, “Quantum embeddings for machine learning,” (2020), arXiv:2001.03622 [quant-ph].
- [51] Phuc Q Le, Fangyan Dong, and Kaoru Hirota, “A flexible representation of quantum images for polynomial preparation, image compression, and processing operations,” *Quantum Information Processing* **10**, 63–84 (2011).
- [52] Yi Zhang, Kai Lu, Yinghui Gao, and Mo Wang, “NEQR: a novel enhanced quantum representation of digital images,” *Quantum Information Processing* **12**, 2833–2860 (2013).
- [53] Harper R. Grimsley, Sophia E. Economou, Edwin Barnes, and Nicholas J. Mayhall, “An adaptive variational algorithm for exact molecular simulations on a quantum computer,” *Nature Communications* **10**, 3007 (2019).
- [54] Sukin Sim, Peter D. Johnson, and Alán Aspuru-Guzik, “Expressibility and entangling capability of parameterized quantum circuits for hybrid quantum-classical algorithms,” *Advanced Quantum Technologies* **2**, 1900070 (2019).
- [55] Maria Schuld, Ville Bergholm, Christian Gogolin, Josh Izaac, and Nathan Killoran, “Evaluating analytic gradients on quantum hardware,” *Phys. Rev. A* **99**, 032331 (2019).
- [56] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner, “Gradient-based learning applied to document recognition,” *Proceedings of the IEEE* **86**, 2278–2324 (1998).
- [57] Alex Krizhevsky, Geoffrey Hinton, *et al.*, “Learning multiple layers of features from tiny images,” (2009).
- [58] Gadi Aleksandrowicz, Thomas Alexander, Panagiotis Barkoutsos, Luciano Bello, Yael Ben-Haim, David Bucher, Francisco Jose Cabrera-Hernández, Jorge Carballo-Franquis, Adrian Chen, Chun-Fu Chen, Jerry M. Chow, Antonio D. Córcoles-Gonzales, Abigail J. Cross, Andrew Cross, Juan Cruz-Benito, Chris Culver, Salvador De La Puente González, Enrique De La Torre, Delton Ding, Eugene Dumitrescu, Ivan Duran, Pieter Eendebak, Mark Everitt, Ismael Faro Sertage, Albert Frisch, Andreas Fuhrer, Jay Gambetta, Borja Godoy Gago, Juan Gomez-Mosquera, Donny Greenberg, Ikko Hamamura, Vojtech Havlicek, Joe Hellmers, Lukasz Herok, Hiroshi Horii, Shaohan Hu, Takashi Imamichi, Toshinari Itoko, Ali Javadi-Abhari, Naoki Kanazawa, Anton Karazeev, Kevin Krsulich, Peng Liu, Yang Luh, Yunho Maeng, Manoel Marques, Francisco Jose Martín-Fernández, Douglas T. McClure, David McKay, Srujan Meesala, Antonio Mezzacapo, Nikolaj Moll, Diego Moreda Rodríguez, Giacomo Nannicini, Paul Nation, Pauline Ollitrault, Lee James O’Riordan, Hanhee Paik, Jesús Pérez, Anna Phan, Marco Pistoia, Viktor Prutyanov, Max Reuter, Julia Rice, Abdón Rodríguez Davila, Raymond Harry Putra Rudy, Mingi Ryu, Ninad Sathaye, Chris Schnabel, Eddie Schoute, Kanav Setia, Yunong Shi, Adenilton Silva, Yukio Siraichi, Seyon Sivarajah, John A. Smolin, Mathias Soeken, Hitomi Takahashi, Ivano Tavernelli, Charles Taylor, Pete Taylor, Kenso Trabing, Matthew Treinish, Wes Turner, Desiree Vogt-Lee, Christophe Vuillot, Jonathan A. Wildstrom, Jessica Wilson, Erick Winston, Christopher Wood, Stephen Wood, Stefan Wörner, Ismail Yunus Akhalwaya, and Christa Zoufal, “Qiskit: An Open-source Framework for Quantum Computing,” (2019).
- [59] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, *et al.*, “Pytorch: An imperative style, high-performance deep learning library,” arXiv preprint arXiv:1912.01703 (2019).
- [60] Alexander Kolesnikov, Xiaohua Zhai, and Lucas Beyer, “Revisiting self-supervised visual representation learning,” in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition* (2019) pp. 1920–1929.
- [61] Ben Jaderberg, Abhishek Agarwal, Karsten Leonhardt, Martin Kiffner, and Dieter Jaksch, “Minimum hardware requirements for hybrid quantum-classical dmft,” *Quantum Science and Technology* **5**, 034015 (2020).
- [62] Frank Arute, Kunal Arya, Ryan Babbush, Dave Bacon, Joseph C Bardin, Rami Barends, Rupak Biswas, Sergio Boixo, Fernando GSL Brandao, David A Buell, *et al.*, “Quantum supremacy using a programmable superconducting processor,” *Nature* **574**, 505–510 (2019).
- [63] Morten Kjaergaard, Mollie E Schwartz, Jochen Braumüller, Philip Krantz, Joel I-J Wang, Simon Gustavsson, and William D Oliver, “Superconducting qubits: Current state of play,” *Annual Review of Condensed Matter Physics* **11**, 369–395 (2020).
- [64] Petar Jurcevic, Ali Javadi-Abhari, Lev S Bishop, Isaac Lauer, Daniela Borgorin, Markus Brink, Lauren Capelluto, Oktay Gunluk, Toshinari Itoko, Naoki Kanazawa, *et al.*, “Demonstration of quantum volume 64 on a superconducting quantum computing system,” *Quantum Science and Technology* (2021).
- [65] Gael Sentís, John Calsamiglia, Ramón Muñoz-Tapia, and Emilio Bagan, “Quantum learning without quantum memory,” *Scientific reports* **2**, 1–8 (2012).
- [66] Unai Alvarez-Rodriguez, Lucas Lamata, Pablo Escandell-Montero, José D Martín-Guerrero, and Enrique Solano, “Supervised quantum learning without measurements,” *Scientific reports* **7**, 1–9 (2017).
- [67] Mohammad H Amin, Evgeny Andriyash, Jason Rolfe, Bohdan Kulchytskyy, and Roger Melko, “Quantum boltzmann machine,” *Physical Review X* **8**, 021050 (2018).
- [68] Gael Sentís, Alex Monras, Ramon Muñoz-Tapia, John Calsamiglia, and Emilio Bagan, “Unsupervised classification of quantum data,” *Physical Review X* **9**, 041029 (2019).
- [69] Iris Cong, Soonwon Choi, and Mikhail D Lukin, “Quantum convolutional neural networks,” *Nature Physics* **15**, 1273–1278 (2019).
- [70] Han-Sen Zhong, Hui Wang, Yu-Hao Deng, Ming-Cheng Chen, Li-Chao Peng, Yi-Han Luo, Jian Qin, Dian Wu, Xing Ding, Yi Hu, *et al.*, “Quantum computational advantage using photons,” *Science* **370**, 1460–1463 (2020).
- [71] Diederik P Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” arXiv preprint arXiv:1412.6980 (2014).

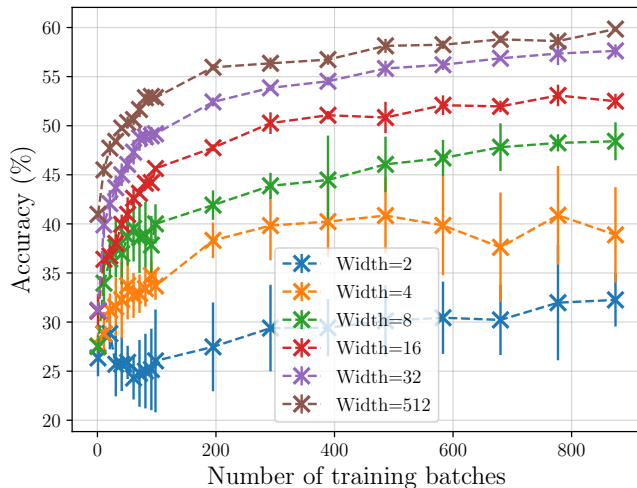


FIG. 7. Classification accuracy achieved in linear probing experiments by classical representation networks with varying network widths at checkpoints across self-supervised training. The markers show the average of three independently trained self-supervised models and linear probe experiments, whilst the error bars show one standard deviation.

### Appendix A: Classical ablation

In order to incorporate QNNs that can be run on current quantum devices into contrastive learning, a compression of the feature vector is required after ConvNet. Since this would not be necessary in a purely classical setting, its impact on final performance is not well understood. To this end, we perform a study of the accuracy achieved by models with different representation network widths. We do this with classical representation networks to remove the quantum specific considerations of statistical noise and optimal circuit architecture, focusing purely on width. The classical representation network is a two-layer, width  $W$  MLP, with Leaky ReLu activation functions after each layer and with bias.

Each model is trained on the first five classes of the CIFAR-10 dataset and a linear probe experiment evaluates the performance at regular checkpoints during training. Fig. 7 shows the result comparing models with different representation network widths, including the  $W = 512$  case which corresponds to no compression. Starting from  $W = 2$ , we see that increasing the width of the representation network improves the test accuracy. Furthermore, we find that  $W = 8$  is the lowest width network in which test accuracy retains the same qualitative behaviour as the uncompressed network. Therefore, in our proof-of-principle quantum experiments, we use an eight width representation network corresponding to

eight qubits.

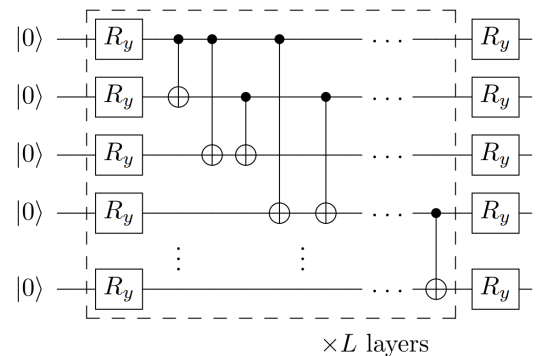


FIG. 8. Alternative variational ansatz. Each layer consists of a single qubit  $\hat{R}_y$  rotation on each qubit, followed by CNOT gates connecting all qubits to each other. After all layers have been applied, a final set of  $\hat{R}_y$  rotations are applied. Every rotation gate is parameterised by a different variational parameter.

### Appendix B: Performance of alternative ansatz

In the main body of this article, all QNNs are constructed using the variational ansatz in Fig. 3, which connects the qubits in a ring of parameterised controlled rotation gates. Here we introduce a second ansatz, as seen in Fig. 8, which is different in that it connects all of the qubits together and only has single qubit parameterised gates. Notably, this ansatz was recently shown to exhibit a larger effective dimension when applied to supervised learning than equivalent classical networks [38]. Therefore, we test whether this circuit structure is also a good candidate for improved performance in a self-supervised setting.

We repeat training, using the same setup as in Fig. 5, with a quantum model containing the new all-to-all ansatz simulated on a statevector simulator. Importantly, for a fair comparison, we apply three layers of the all-to-all ansatz so that both circuits have 32 learnable parameters. The result of the linear probe experiments can be seen in Fig. 9, along with the previous models for comparison. We see that for the all-to-all ansatz, test accuracy is no higher than the classical model beyond the statistical variance of repeating training with different initial parameters, and below the ring ansatz. Indeed, by the end of training, the all-to-all ansatz achieves a final accuracy of  $(43.46 \pm 1.68)\%$ , which is similar to the classical model. Thus, we show that achieving an advantage using quantum neural networks in contrastive learning is highly dependent on the correct choice of quantum circuit structure.

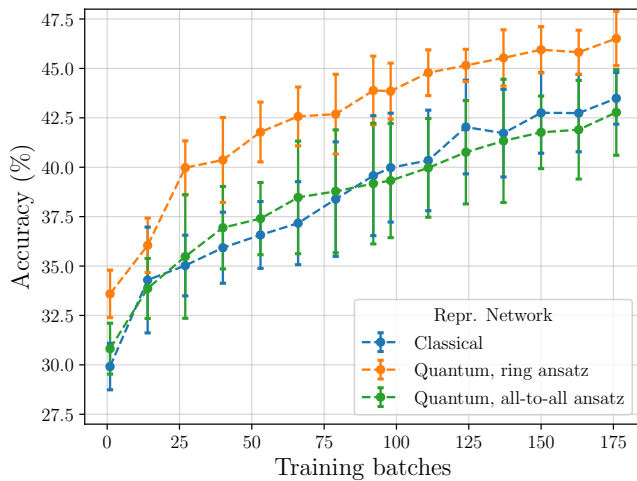


FIG. 9. Classification accuracy achieved in linear probing experiments using the encoder at checkpoints across self-supervised training. Comparison between models trained with a classical representation network (blue), quantum representation network with the ring ansatz (orange) and quantum representation network with the all-to-all ansatz (green). All quantum circuits were evaluated on a statevector simulator. The markers show the average of six independently trained models, whilst the error bars show one standard deviation.