# Shadows of quantum machine learning

Jerbi, S.; Gyurik, C.F.S.; Marshall, S.C.; Molteni, R.; Dunjko, V.

# Shadows of quantum machine learning

Sofiene Jerbi [1,2] ✉, Casper Gyurik[3], Simon C. Marshall[3], Riccardo Molteni[3] & Vedran Dunjko [3]

Quantum machine learning is often highlighted as one of the most promising practical applications for which quantum computers could provide a computational advantage. However, a major obstacle to the widespread use of quantum machine learning models in practice is that these models, even once trained, still require access to a quantum computer in order to be evaluated on new data. To solve this issue, we introduce a class of quantum models where quantum resources are only required during training, while the deployment of the trained model is classical. Specifically, the training phase of our models ends with the generation of a 'shadow model' from which the classical deployment becomes possible. We prove that: (i) this class of models is universal for classically-deployed quantum machine learning; (ii) it does have restricted learning capacities compared to 'fully quantum' models, but nonetheless (iii) it achieves a provable learning advantage over fully classical learners, contingent on widely believed assumptions in complexity theory. These results provide compelling evidence that quantum machine learning can confer learning advantages across a substantially broader range of scenarios, where quantum computers are exclusively employed during the training phase. By enabling classical deployment, our approach facilitates the implementation of quantum machine learning models in various practical contexts.

Quantum machine learning is a rapidly growing field[1–3] driven by its potential to achieve quantum advantages in practical applications. A particularly interesting approach to make quantum machine learning applicable in the near term is to develop learning models based on parametrized quantum circuits[4–6]. Indeed, such quantum models have already been shown to achieve good learning performance in benchmarking tasks, both in numerical simulations[7–11] and on actual quantum hardware[12–15]. Moreover, based on widely believed cryptography assumptions, these models also hold the promise to solve certain learning tasks that are intractable for classical algorithms[16,17], including predicting ground state properties of highly-interacting quantum systems[18].

Despite these advances, quantum machine learning is facing a major obstacle for its use in practice. A typical workflow of a machine learning model involved, e.g., in driving autonomous vehicles, is divided into: (i) a training phase, where the model is trained, typically using training data or by reinforcement; followed by (ii) a deployment phase,

where the trained model is evaluated on new input data. For quantum machine learning models, both of these phases require access to a quantum computer. But given that in many practical machine learning applications, the trained model is meant for a widespread deployment, the current scarcity of quantum computing access dramatically reduces the applicability of quantum machine learning. One way of addressing this problem is by generating shadow models out of quantum machine learning models. That is, we propose inserting a shadowing phase between the training and deployment, where a quantum computer is used to collect information on the quantum model. Then a classical computer can use this information to evaluate the model on new data during the deployment phase.

The conceptual idea of generating shadows of quantum models was already proposed by Schreiber et al.[19], albeit under the terminology of classical surrogates. In that work, as well as in that of Landman et al.[20], the authors make use of the general expression of quantum models as trigonometric polynomials[21] to learn the Fourier

[1]Institute for Theoretical Physics, University of Innsbruck, Innsbruck, Austria. [2]Dahlem Center for Complex Quantum Systems, Freie Universität Berlin, Berlin, Germany. [3]applied Quantum algorithms (aQa), Leiden University, Leiden, The Netherlands. ✉e-mail: sofiene.jerbi@fu-berlin.de
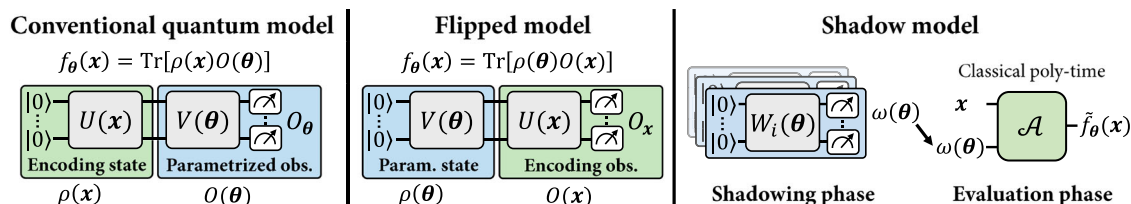
**Fig. 1 | Quantum and shadow models.** (left) Conventional quantum models can be expressed as inner products between a data-encoding quantum state $\rho(x)$ and a parametrized observable $O(\theta)$. The resulting linear model $f_\theta(x) = \mathrm{Tr}[\rho(x)O(\theta)]$ naturally corresponds to a quantum computation, depicted here. (middle) We define flipped models $f_\theta(x) = \mathrm{Tr}[\rho(\theta)O(x)]$ as quantum linear models where the role of the quantum state $\rho(\theta)$ and the observable $O(x)$ is flipped compared to conventional models. (right) Flipped models are associated to natural shadow models: one can use techniques from shadow tomography to construct a classical representation $\hat\rho(\theta)$ of the parametrized state $\rho(\theta)$ (during the shadowing phase), such

that, for encoding observables $O(x)$ that are classically representable (e.g., linear combinations of Pauli observables), $\hat\rho(\theta)$ can be used by a classical algorithm to evaluate the model $f_\theta(x)$ on new input data (during the evaluation phase). More generally, a shadow model is defined by (i) a shadowing phase where a (bit-string) advice $\omega(\theta)$ is generated by the evaluation of multiple quantum circuits $W_1(\theta), \ldots, W_M(\theta)$, and (ii) an evaluation phase where this advice is used by a classical algorithm $\mathcal{A}$, along with new input data $x$ to evaluate their labels $\tilde f_\theta(x)$. In the Section "General shadow models", we show that under this general definition, all shadow models are shadows of flipped models.

representation of trained models and evaluate them classically on new data. However, these works also suggest that a classical model could potentially be trained directly on the training data and achieve the same performance as the shadow model, thus circumventing the need for a quantum model in the first place. This raises the concern that all quantum models that are compatible with a classical deployment would also lose all quantum advantage, hence severely limiting the prospects for a widespread use of quantum machine learning.

Therefore, two natural open questions are raised:

1. Can shadow models achieve a quantum advantage over entirely classical (classically trained and classically evaluated) models?
2. Do there exist quantum models that do not admit efficiently evaluatable shadow models?

In this work, we resolve both of these key open questions. We propose a general definition for shadow models, rooted in the fundamental idea that quantum machine learning models can be universally expressed as linear models[22]. This formulation of shadow models allows us to leverage various results and techniques from quantum information theory for the analysis of this model class. From a practical perspective, employing shadow tomography techniques[23–26] allows to easily construct diverse shadow models that will resonate with the practitioners of quantum machine learning. Furthermore, in our exploration of the computational capabilities of shadow models, we find them to capture a distinct computational class. Specifically, we demonstrate that, under widely believed cryptography assumptions, there exist learning tasks where shadow models exhibit a provable quantum advantage over fully classical models. However, contrary to this advantage, we also establish that there exist quantum models that are strictly more powerful than the class of shadow models, based on common assumptions in complexity theory.

For ease of exposition, we will first adhere to a working definition of a shadow model as a model that is trained on a quantum computer, but can be evaluated classically on new input data with the help of information generated by a quantum computer (i.e., quantum-generated advice) that is independent of the new data. We will (informally) call a model "shadowfiable" if there exists a method of turning it into a shadow model. In the Section "General shadow models", we will make our definitions more precise.

## Results

### The flipped model

The construction of our shadow models starts from a simple yet key observation: all standard quantum machine learning models for supervised learning can be expressed as linear models[22]. To delve into this claim, we first draw upon early works that utilized parametrized quantum circuits in machine learning[7,12]. These works proposed

quantum models that are naturally expressed as linear functions of the form

$$f_\theta(x) = \mathrm{Tr}[\rho(x)O(\theta)] \tag{1}$$

where $\rho(x)$ are quantum states that encode classical data $x \in \mathcal{X}$ and $O(\theta)$ are parametrized observables whose inner product with $\rho(x)$ defines $f_\theta(x)$ (see Fig. 1). In a regression task, one would use such a model to assign a real-valued label to an input $x$, while in classification tasks, one would additionally apply, e.g., a sign function, to discretize its output into a class. From a circuit picture, such models can be evaluated on a quantum computer by: (i) preparing an initial state $\rho_0$, e.g., $|0\rangle\langle 0|^{\otimes n}$, (ii) evolving it under a data-dependent circuit $U(x)$, (iii) followed by a variational circuit $V(\theta)$, (iv) before finally measuring the expectation value of a Hermitian observable $O$. Together, steps (i) and (ii) define

$$\rho(x) = U(x)\rho_0 U^\dagger(x), \tag{2}$$

while steps (iii) and (iv) define

$$O(\theta) = V^\dagger(\theta)OV(\theta). \tag{3}$$

Since the early works, it is known that quantum linear models also capture quantum kernel models as a special case[27], simply by making $O(\theta)$ directly dependent on the training data of the learning task. Perhaps more surprisingly, quantum linear models can also encompass more general data re-uploading models, composed of several layers of data encoding and variational processing $U_1(x)V_1(\theta)U_2(x)\ldots$ Indeed, data re-uploading models can be mapped to linear models through circuit transformations (e.g., gate teleportation) that relocate all data-encoding gates to the first layer of the circuit[22].

### Flipped model definition

The definition of a quantum linear model in Eq. (1) can in general accommodate any pair of Hermitian operators in place of $\rho(x), O(\theta)$. However, due to how these models are evaluated on a quantum computer, one commonly works under the constraint that $\rho(x)$ defines a quantum state (i.e., a positive semi-definite operator with unit trace). Indeed, from an operational perspective, $\rho(x)$ must be physically prepared on a quantum device before being measured with respect to the observable $O(\theta)$ (which only needs to be Hermitian in order to be a valid observable).

For reasons that will become clearer from the shadowing perspective, we define a so-called flipped model, where we flip the role of

$\rho(\boldsymbol{x})$ and $O(\boldsymbol{\theta})$. That is, we consider

$$f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \mathrm{Tr}[\rho(\boldsymbol{\theta})O(\boldsymbol{x})] \tag{4}$$

where $\rho(\boldsymbol{\theta})$ is a parametrized quantum state and $O(\boldsymbol{x})$ is an observable that encodes the data and can take more general forms than Eq. (3) as we will see next. This model also corresponds to a straightforward quantum computation as $\rho(\boldsymbol{\theta})$ can be physically prepared before being measured with respect to $O(\boldsymbol{x})$.

A simple example of flipped model is for instance defined by:

$$\rho(\boldsymbol{\theta}) = V(\boldsymbol{\theta})\rho_0 V^{\dagger}(\boldsymbol{\theta}) \quad \& \quad O(\boldsymbol{x}) = \sum_{j=1}^{m} w_j(\boldsymbol{x})P_j \tag{5}$$

for an initial state $\rho_0$, a variational circuit $V(\boldsymbol{\theta})$, and a collection of Pauli observables $\{P_j\}_{j=1}^{m}$ weighted by data-dependent weights $w_j(\boldsymbol{x}) \in \mathbb{R}$. One can evaluate this model by repeatedly preparing $\rho(\boldsymbol{\theta})$ on a quantum computer, measuring it in a Pauli basis specified by a $P_j$, and weighting the outcome by $w_j(\boldsymbol{x})$. For other examples of flipped models, see Supplementary Section 1.

As opposed to conventional quantum linear models, flipped models are well-suited to construct shadow models. Since the variational operators $\rho(\boldsymbol{\theta})$ are quantum states, one can straightforwardly use techniques from shadow tomography[23] to construct classical shadows $\hat{\rho}(\boldsymbol{\theta})$ of these states. What we call classical shadows $\hat{\rho}(\boldsymbol{\theta})$ here are collections of measurement outcomes obtained from copies of $\rho(\boldsymbol{\theta})$ that can be used to classically approximate expectation values of certain observables $O$ (for a certain restricted family). If we take these observables to be our data-dependent $O(\boldsymbol{x})$, then we end up with a classical model $\widetilde{f}_{\boldsymbol{\theta}}(\boldsymbol{x})$ that approximates our flipped model. Note here that one has total freedom on the classical shadow techniques they may use to define their shadow models, and a plethora of protocols have already been proposed in the literature[23–26]. But it is important to keep in mind that each of these protocols comes with its limitations, as it may restrict the class of states $\rho(\boldsymbol{\theta})$ or the class of observables $O(\boldsymbol{x})$ for which an efficient and faithful shadow model can be constructed. By "efficient" we refer here to the number of measurements performed on $\rho(\boldsymbol{\theta})$ and the time complexity of estimating the expectation values of observables $O(\boldsymbol{x})$ from these measurements. And by "faithful" we refer to the approximation error between the shadow model $\widetilde{f}_{\boldsymbol{\theta}}(\boldsymbol{x})$ resulting from the shadow protocol and the original flipped model $f_{\boldsymbol{\theta}}(\boldsymbol{x})$. For instance, in the example of Eq. (5), we know that if all Pauli operators $\{P_j\}_{j=1}^{m}$ are $k$-local, then $\widetilde{\mathcal{O}}(3^k B^2 \varepsilon^{-2})$ measurements of $\rho(\boldsymbol{\theta})$, where $B = \max_{\boldsymbol{x}} \sum_{j=1}^{m} |w_i(\boldsymbol{x})|$, are sufficient to guarantee $\max_{\boldsymbol{x}} |\widetilde{f}_{\boldsymbol{\theta}}(\boldsymbol{x}) - f_{\boldsymbol{\theta}}(\boldsymbol{x})| \leq \varepsilon$ with high probability. But for non-local Pauli operators (i.e., large $k$), this protocol becomes highly inefficient if we want to guarantee a low error $\varepsilon$.

Importantly, shadowfied flipped models are not limited to constructions based on classical shadow protocols. Given that the states $\rho(\boldsymbol{\theta})$ are not given to us a black-box (as is generally assumed in shadow tomography), one can use prior knowledge on these states to construct efficient shadowing procedure. For instance, if $\rho(\boldsymbol{\theta})$ is known to be a superposition of a tractable number of computational basis states, or well-approximated by a matrix product state (MPS) with low bond dimension, then efficient tomography protocols may be used[28].

## Properties of flipped models

Flipped models are a stepping stone toward the claims of quantum advantage and "shadowfiability" that are the focus of this paper. Nonetheless, they constitute a newly introduced model, which is why it is useful to understand first how they relate to previous quantum models and what learning guarantees they can have.

Since conventional linear models of the form of Eq. (1) play a central role in quantum machine learning, we start by asking the question: when can these models be represented by (efficiently evaluatable) flipped models? That is, given a conventional model $f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \mathrm{Tr}[\rho(\boldsymbol{x})O(\boldsymbol{\theta})]$, can we construct a flipped model $\widetilde{f}_{\boldsymbol{\theta}}(\boldsymbol{x}) = \mathrm{Tr}[\rho'(\boldsymbol{\theta})O'(\boldsymbol{x})]$ such that $\widetilde{f}_{\boldsymbol{\theta}}(\boldsymbol{x}) \approx f_{\boldsymbol{\theta}}(\boldsymbol{x}), \forall \boldsymbol{x}, \boldsymbol{\theta}$, and $\widetilde{f}_{\boldsymbol{\theta}}(\boldsymbol{x})$ is as efficient to evaluate as $f_{\boldsymbol{\theta}}(\boldsymbol{x})$? Clearly, a conventional model $f_{\boldsymbol{\theta}}(\boldsymbol{x})$ for which the parametrized operator $O(\boldsymbol{\theta})$ is also a quantum state (i.e., a positive semi-definite trace-1 operator) is by definition also a flipped model. Therefore, a natural strategy to flip a conventional model is to transform its observable $O(\boldsymbol{\theta})$ into a quantum state $\rho'(\boldsymbol{\theta})$. This transformation involves dealing with the negative eigenvalues of $O(\boldsymbol{\theta})$, which can be taken into account using an auxiliary qubit, without overheads in the efficiency of evaluation (see Supplementary Section 2 for more details). More importantly, the transformation involves normalizing these eigenvalues, which affects the efficiency of evaluating the resulting flipped model. Indeed, the normalization factor $\alpha$ that results from normalizing $O(\boldsymbol{\theta})$ corresponds to its trace norm $\| O \|_1 = \mathrm{Tr}\left[\sqrt{O^2}\right]$ and needs to be absorbed into the observable $O'(\boldsymbol{x}) = \alpha \rho(\boldsymbol{x})$ of the flipped model $\widetilde{f}_{\boldsymbol{\theta}}(\boldsymbol{x})$ to guarantee $\widetilde{f}_{\boldsymbol{\theta}}(\boldsymbol{x}) = f_{\boldsymbol{\theta}}(\boldsymbol{x})$. This directly impacts the spectral norm $\| O' \|_{\infty} = \max_{|\psi\rangle} \langle O' \rangle_{\psi} = \alpha$ of the flipped model, and therefore the efficiency of its evaluation, as $\mathcal{O}(\| O' \|_{\infty}^2 / \varepsilon^2)$ measurements of $\rho'(\boldsymbol{\theta})$ are needed in order to estimate $\widetilde{f}_{\boldsymbol{\theta}}(\boldsymbol{x})$ to additive error $\varepsilon$ (see Supplementary Section 2 for a derivation). Therefore, we end up showing that, for a conventional model $f_{\boldsymbol{\theta}}(\boldsymbol{x})$ acting on $n$ qubits and with a bounded observable trace norm $\|O\|_1 \leq \alpha$, we can construct a flipped model acting on $m = n + 1$ qubits and with observable spectral norm $\| O' \|_{\infty} = \alpha$.

Interestingly, in the relevant regime where the number of qubits $n, m$ used by the linear models involved in this flipping is logarithmic in $\|O\|_1$ (e.g., where $O$ is a Pauli observable and hence $\|O\|_1 = 2^n$), we find that this requirement on the spectral norm $\| O' \|_{\infty}$ of the resulting flipped model is unavoidable in the worst case, up to a logarithmic factor in $\|O\|_1$. We refer to Appendix Supplementary Section 2 for proofs of these statements and a more in-depth discussion.

Another property of interest in machine learning is the generalization performance of a learning model. That is, we want to bound the gap between the performance of the model on its training set (so-called training error) and its performance on the rest of the data space (or expected error). Such bounds have for instance been derived in terms of the number of encoding gates in the quantum model[29], or the rank of its observable[30]. In the case of flipped model, we find instead a bound in terms of the number of qubits $n$ and the spectral norm $\|O\|_{\infty}$ of the observable. Since these quantities are operationally meaningful, this gives us a natural way of controlling the generalization performance of our flipped models. Stated informally, we find that if a flipped model achieves a small error $|f_{\boldsymbol{\theta}}(\boldsymbol{x}) - f(\boldsymbol{x})| \leq \eta$ for all $\boldsymbol{x}$ in a training set of size $M$, then we only need $M$ to scale as $\widetilde{\Omega}\left(\frac{n\|O\|_{\infty}^2}{\varepsilon\eta^2}\right)$ in order to guarantee a small expected error $|f_{\boldsymbol{\theta}}(\boldsymbol{x}) - f(\boldsymbol{x})| \leq 2\eta$ with probability $1 - \varepsilon$ over the entire data distribution.

Note that the dependence on $n$ and $\|O\|_{\infty}$ is linear and quadratic, respectively, which means that we can afford a large number of qubits and a large spectral norm and still guarantee a good generalization performance. This is particularly relevant as the spectral norm is a controllable quantity, meaning we can easily fine-tune our models to perform well in training and generalize well. E.g., in the case of the model in Eq. (5), this spectral norm is bounded by $\max_{\boldsymbol{x}} \sum_{j=1}^{m} |w_i(\boldsymbol{x})|$, which scales favourably with the number of qubits $n$ if $m \in \mathcal{O}(\mathrm{poly}(n))$ or if the vector $\boldsymbol{w}(\boldsymbol{x})$ is sparse.

## Quantum advantage of a shadow model

We recall that we (informally) define shadow models as models that are trained on a quantum computer, but, after a shadowing

procedure that collects information on the trained model, are evaluated classically on new input data. In this section, we consider the question of achieving a quantum advantage using such shadow models. It may seem at first sight that this question has a straightforward answer, which is "no": if the function learned by a model is classically computable, then there should be no room for a quantum advantage. However, as demonstrated in refs. [17,31], one can also achieve a quantum advantage based on so-called trap-door functions. These are functions that are believed to be hard to compute classically, unless given a key (or advice) that allows for an efficient classical computation. Notably, there exist trap-door functions where this key can be efficiently computed using a quantum computer, but not classically. This allows us to construct shadow models that make use of this quantum-generated key to compute an otherwise classically untractable function.

Similarly to related results showing a quantum advantage in machine learning with classical data[16,32], we consider a learning task where the target function (i.e., the function generating the training data) is derived from cryptographic functions that are widely believed to be hard to compute classically. More precisely, we introduce a variant of the discrete cube root learning task[17], which is hard to solve classically under a hardness assumption related to that of the RSA cryptosystem[33]. In this task, we consider target functions defined on $\mathbb{Z}_N = \{0, \ldots, N-1\}$ as

$$g_s(\boldsymbol{x}) = \begin{cases} 1, & \text{if } \sqrt[3]{\boldsymbol{x}} \bmod N \in [s, s + \frac{N-1}{2}], \\ 0, & \text{otherwise} \end{cases} \tag{6}$$

where $N = pq$ is an $n$-bit integer, product of two primes $p, q$ of the form $3k+2, 3k'+2$, such that the discrete cube root is properly defined as the inverse of the function $\boldsymbol{y}^3 \bmod N$. These target functions are particularly appealing because of a number of interesting properties:

i.  It is believed that given only $\boldsymbol{x}$ and $N$ as input, computing $g(\boldsymbol{x}) = \sqrt[3]{\boldsymbol{x}} \bmod N$ with high probability of success over random draws of $\boldsymbol{x}$ and $N$ is classically intractable. This assumption is known as the discrete cube root (DCR) assumption.

ii. On the other hand, computing $\boldsymbol{x}^a \bmod N$ is classically efficient for any $a \in \mathbb{Z}_N$. For $a = 3$, this implies that $g^{-1}(\boldsymbol{y}) = \boldsymbol{y}^3 \bmod N$ is a one-way function, under the DCR assumption.

iii. The function $g(\boldsymbol{x}) = \sqrt[3]{\boldsymbol{x}} \bmod N$ has a "trap-door", in that there exists another way of computing it efficiently. For every $N$ (as specified above), there exists a key $d \in \mathbb{Z}_N$ such that $g(\boldsymbol{x}) = \boldsymbol{x}^d \bmod N$. Finding $d$ is efficient quantumly by using Shor's factoring algorithm[34], but hard classically under the DCR assumption.

Observations (i) and (ii) can be leveraged to show that learning the functions $g_s$ from examples is also intractable. Indeed, Alexi et al.[35] showed that a classical algorithm that could faithfully capture a single bit $g_s(\boldsymbol{x})$ of the discrete cube root of $\boldsymbol{x}$, for even a $1/2 + 1/\text{poly}(n)$ fraction of all $\boldsymbol{x} \in \mathbb{Z}_N$, could also be used to reconstruct $g(\boldsymbol{x}), \forall \boldsymbol{x} \in \mathbb{Z}_N$, with high probability of success. Since, from observation (ii), the training data for the learning algorithm can also be generated efficiently classically from $N$, a classical learner that learns $g_s(\boldsymbol{x})$ correctly for a $1/2 + 1/\text{poly}(n)$ fraction of all $\boldsymbol{x} \in \mathbb{Z}_N$ would then contradict the DCR assumption.

Observation (iii) allows us to define the following flipped model:

$$f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \text{Tr}[\rho(\boldsymbol{\theta})O(\boldsymbol{x})]$$
$$\rho(\boldsymbol{\theta}) = |d', s'\rangle\langle d', s'| \ \& \ O(\boldsymbol{x}) = \sum_{d', s'} \widehat{g}_{d', s'}(\boldsymbol{x})|d', s'\rangle\langle d', s'|. \tag{7}$$

That is, $\rho(\boldsymbol{\theta})$ (for $\boldsymbol{\theta} = (N, s')$) specifies candidates for the key $d'$ and the parameter $s'$ of interest, while $O(\boldsymbol{x})$ uses that information to compute

$$\widehat{g}_{d', s'}(\boldsymbol{x}) = \begin{cases} 1, & \text{if } \boldsymbol{x}^{d'} \bmod N \in [s', s' + \frac{N-1}{2}], \\ 0, & \text{otherwise.} \end{cases} \tag{8}$$

The state $\rho(\boldsymbol{\theta}) = |d, s'\rangle\langle d, s'|$ for the right key $d$ can be prepared efficiently using Shor's algorithm applied on $N$ (provided with the training data). As for $O(\boldsymbol{x})$, it simply processes classically a bit-string to compute $\widehat{g}_{d', s'}(\boldsymbol{x})$ efficiently, which corresponds to $g_s(\boldsymbol{x})$ when $(d', s') = (d, s)$. Finding an $s'$ close to $s$ is an easy task given training data and $d' = d$. Since $\rho(\boldsymbol{\theta})$ is a computational basis state, this flipped model admits a trivial shadow model where a single computational basis measurement of $\rho(\boldsymbol{\theta})$ allows to evaluate $f_{\boldsymbol{\theta}}(\boldsymbol{x})$ classically for all $\boldsymbol{x}$. Therefore, we end up showing the following theorem:

**Theorem 1. (Quantum advantage (informal)).** There exists a learning task where a shadow model first trained using a quantum computer then evaluated classically on new input data, can achieve an arbitrarily good learning performance, while any fully classical model cannot do significantly better than random guessing, under the hardness of classically computing the discrete cube root.

In Supplementary Section 3, we formalize the statement of this result using the PAC framework and provide more details on the setting and the proofs.

## General shadow models

As mentioned at the start of this paper, shadow models are not limited to shadowfied flipped models, and the main alternative proposals are based on the Fourier representation of quantum models[19,20]. It is clear that Fourier models are defined very differently from flipped models, but one may wonder whether they nonetheless include shadowfied flipped models as a special case, or the other way around.

In this section, we first start by showing that there exist quantum models that admit shadow models (i.e., are shadowfiable) but cannot be shadowfied efficiently using a Fourier approach. This then motivates our proposal for a general definition of shadow models, and we show that, under this definition, all shadow models can be expressed as shadowfied flipped models. Finally, we show the existence of quantum models that are not shadowfiable at all under likely complexity theory assumptions.

### Shadow models beyond Fourier

An interesting approach to construct shadows of quantum models is based on their natural Fourier representation. It has been shown[21,36] that quantum models can be expressed as generalized Fourier series of the form

$$f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \sum_{\boldsymbol{\omega} \in \Omega} c_{\boldsymbol{\omega}}(\boldsymbol{\theta})e^{-i\boldsymbol{\omega} \cdot \boldsymbol{x}} \tag{9}$$

where the accessible frequencies $\Omega$ only depend on properties of the encoding gates used by the model (notably the number of encoding gates and their eigenvalues). Since these frequencies can easily be read out from the circuit, one can proceed to form a shadow model by estimating their associated coefficients $c_{\boldsymbol{\omega}}(\boldsymbol{\theta})$ using queries of the quantum model $f_{\boldsymbol{\theta}}(\boldsymbol{x})$ at different values $\boldsymbol{x}$ and, e.g., a Fourier transform[19]. Given a good approximation of these coefficients, one can then compute estimates of $f_{\boldsymbol{\theta}}(\boldsymbol{x})$ for arbitrary new inputs $\boldsymbol{x}$. We will refer to such a shadowing approach that considers the quantum model as a black-box, aside from the knowledge of its Fourier spectrum, as the Fourier shadowing approach.

Although we will be explicit about this in the next subsection, we will consider a shadowing procedure to be successful, if, with high probability, the resulting shadow model agrees with the original model on all inputs, i.e.,

$$\max_{\boldsymbol{x}\in\mathcal{X}}|f_{\boldsymbol{\theta}}(\boldsymbol{x}) - \widetilde{f}_{\boldsymbol{\theta}}(\boldsymbol{x})| \le \varepsilon, \tag{10}$$

for a specified $\varepsilon \ge 0$. We want the shadowing procedure to be successful independently of the data distribution under which the model should be trained, which justifies this definition. We discuss this point further in Supplementary Section 4.

We show that the Fourier shadowing approach can suffer from an exponential sample complexity in the dimension of the input data $\boldsymbol{x}$, making it intractable for high-dimensional input spaces. To see this, consider the linear model:

$$\begin{aligned} f_{\boldsymbol{y}}(\boldsymbol{x}) &= \mathrm{Tr}[\rho(\boldsymbol{x})O(\boldsymbol{y})] \\ \rho(\boldsymbol{x}) &= \bigotimes_{i=1}^{n} R_Y(x_i)|0\rangle\langle0|R_Y^{\dagger}(x_i) \,\&\, O(\boldsymbol{y}) = |\boldsymbol{y}\rangle\langle\boldsymbol{y}|. \end{aligned} \tag{11}$$

for $\boldsymbol{x} \in \mathbb{R}^n$ and $\boldsymbol{y} \in \{0,1\}^{\otimes n}$. Let us first restrict our attention to the domain $\boldsymbol{x} \in \{0,\pi\}^n$. It is quite clear that on this domain, $f_{\boldsymbol{y}}(\boldsymbol{x}) = \delta_{\boldsymbol{x}/\pi, \boldsymbol{y}}$ plays the role of a database search oracle, where the database has $2^n$ elements and a unique marked element $\boldsymbol{y}$. From lower bounds on database search, we know that $\Omega(2^n)$ calls to this oracle are needed to find $\boldsymbol{y}$[37]. This implies that a Fourier shadowing approach would require $\Omega(2^n)$ calls to $f_{\boldsymbol{y}}(\boldsymbol{x}) = \delta_{\boldsymbol{x}/\pi, \boldsymbol{y}}$ in order to guarantee $\max_{\boldsymbol{x}\in\mathcal{X}}|\widetilde{f}_{\boldsymbol{\theta}}(\boldsymbol{x}) - f_{\boldsymbol{\theta}}(\boldsymbol{x})| \le 1/4$. In Supplementary Section 4, we explain how this result can be generalized to the full domain $\boldsymbol{x} \in \mathbb{R}^n$, and we relate this bound on the sample complexity to the Fourier decomposition of the model.

On the other hand, note that the flipped model associated to $f_{\boldsymbol{y}}(\boldsymbol{x})$ allows for a straightforward shadowing procedure. Indeed, by preparing $O(\boldsymbol{y})$ and measuring it in the computational basis, one straightforwardly obtains $\boldsymbol{y}$ and can therefore classically compute the expectation value of any tensor product observable $\rho(\boldsymbol{x})$ as specified by Eq. (11). Therefore, we have shown that there exist shadowfiable models that are not efficiently Fourier-shadowfiable, i.e., for which a shadowing procedure based solely on the knowledge of their Fourier spectrum and on black-box queries has query complexity that is exponential in the input dimension.

## All shadow models are shadows of flipped models

We give a general definition of shadow models that can encompass all methods that have been proposed to generate them. In contrast to the definition of classical surrogates proposed by Schreiber et al.[19], we give explicit definitions for the shadowing and evaluation phases of shadow models which makes explicit the need for a quantum computer in the shadowing phase. Indeed, as mentioned in the introduction, the term "classical surrogate" has been used to describe both a classically evaluatable model obtained from a quantum shadowing procedure and a fully classical model trained directly on the data. We want to avoid this confusion in the definition of shadow models. We view a general shadowing phase as the generation of advice that can be used to classically evaluate a quantum model. This advice is generated by the execution of quantum circuits that may or may not depend on the (trained) quantum circuit from the training phase. For instance, when we shadowfy a flipped model, we simply prepare the parametrized states $\rho(\boldsymbol{\theta})$ and use (randomized) measurements to generate an operationally meaningful classical description. In the case of Fourier shadowing, this advice is instead generated by evaluations of the quantum model $f_{\boldsymbol{\theta}}(\boldsymbol{x})$ for different inputs $\boldsymbol{x} \in \mathbb{R}^d$ that are rich enough to learn the Fourier coefficients of this model. We propose the following definition:

**Definition 2. (General shadow model).** Let $W_1(\boldsymbol{\theta}), \dots, W_M(\boldsymbol{\theta})$ be a sequence of $\mathcal{O}(\mathrm{poly}\,(m))$-time quantum circuits applied on all-zero

states $|0\rangle^{\otimes m}$, and that can potentially be chosen adaptively. Call $\omega(\boldsymbol{\theta}) = (\omega_1(\boldsymbol{\theta}), \dots, \omega_M(\boldsymbol{\theta}))$ the outcomes of measuring the output states of these circuits in the computational basis. A general shadow model is defined as:

$$f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \mathcal{A}(\boldsymbol{x}, \omega(\boldsymbol{\theta})) \tag{12}$$

where $\mathcal{A}$ is a classical $\mathcal{O}(\mathrm{poly}\,(M, m, d))$ time algorithm that processes the outcomes $\omega(\boldsymbol{\theta})$ along with an input $\boldsymbol{x} \in \mathbb{R}^d$ to return the (real-valued) label $f_{\boldsymbol{\theta}}(\boldsymbol{x})$.

From this definition, a shadow model is a classically evaluatable model that uses quantum-generated advice. Crucially, this advice must be independent of the data points $\boldsymbol{x}$ we wish to evaluate the model on in the future. We distinguish the notion of a shadow model from that of a shadowfiable quantum model, that is a quantum model that admits a shadow model:

**Definition 3. (Shadowfiable model).** A model $f_{\boldsymbol{\theta}}$ acting on $n$ qubits is said to be shadowfiable if, for $\varepsilon, \delta > 0$, there exists a shadow model $\widetilde{f}_{\boldsymbol{\theta}}$ such that, with probability $1 - \delta$ over the quantum generation of the advice $\omega(\boldsymbol{\theta})$ (i.e., the shadowing phase), the shadow model satisfies?

$$\max_{\boldsymbol{x}\in\mathcal{X}}|f_{\boldsymbol{\theta}}(\boldsymbol{x})| - \widetilde{f}_{\boldsymbol{\theta}}(\boldsymbol{x}) \le \varepsilon, \tag{13}$$

and uses $m, M \in \mathcal{O}(\mathrm{poly}\,(n, 1/\varepsilon, 1/\delta))$ qubits and circuits to generate its advice $\omega(\boldsymbol{\theta})$.

While we have seen that there exist shadowfiable models that cannot be shadowfied efficiently using a Fourier approach, we show that all shadowfiable models as defined above can be approximated by shadowfiable flipped models.

**Lemma 4. (Flipped models are shadow-universal).** All shadowfiable models as defined in Defs. 2 and 3 can be approximated by flipped models $f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \mathrm{Tr}[\rho(\boldsymbol{\theta})O(\boldsymbol{x})]$ with the guarantee that computational basis measurements of $\rho(\boldsymbol{\theta})$ and efficient classical post-processing can be used to evaluate $f_{\boldsymbol{\theta}}(\boldsymbol{x})$ to good precision with high probability.

This result is essentially based on the observation that the evaluation of a general shadow model as defined in Def. 2 can be done entirely coherently. Instead of classically running the algorithm $\mathcal{A}$ using the random advice $\omega(\boldsymbol{\theta})$, one can quantumly simulate this algorithm (using a reversible execution) and execute it on the coherent advice $\rho(\boldsymbol{\theta}) = |\omega(\boldsymbol{\theta})\rangle\langle\omega(\boldsymbol{\theta})|$ generated by $\{W_1(\boldsymbol{\theta}), \dots, W_M(\boldsymbol{\theta})\}$ before the computational basis measurements. We refer to Supplementary Section 4 for a more detailed statement and proof.

## Not all quantum models are shadowfiable

From the discrete cube root learning task, we already understand that a learning separation can be established between classical and shadowfiable models. We would also like to understand whether a learning separation exists between shadowfiable models and general quantum models, or equivalently, whether all quantum models are shadowfiable. We show that this also is not the case, under widely believed assumptions (see Fig. 2).

**Theorem 5. (Not all shadowfiable).** Under the assumption that BQP $\not\subset$ P/poly, there exist quantum models, i.e., models in BQP, that are not shadowfiable, i.e., that are not in BPP/qgenpoly.

We start by noting that shadow models can be characterized by a complexity class we define as BPP/qgenpoly, which stands for "Bounded-error Probabilistic Polynomial-time with quantumly generated (polynomial-time) advice of polynomial size". This class contains all functions that can be computed efficiently classically with the help of polynomially-sized advice generated efficiently by a quantum computer. This class is trivially contained in the standard class BPP/poly, which doesn't have any constraint on how the advice is generated and
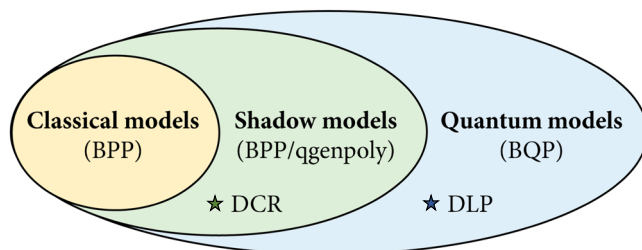
**Fig. 2 | Separations between classical, shadow, and quantum models.** Under the assumption that the discrete cube root (DCR) cannot be computed classically in polynomial time, we have a separation between shadow models (captured by the class BPP/qgenpoly) and classical models (in BPP). Under the assumption that there exist functions that can be computed in quantum polynomial time but not in classical polynomial time with the help of advice (i.e., BQP ⊄ P/poly), we have a separation between quantum models (universal for BQP) and shadow models (BPP/qgenpoly). A candidate function for this separation is the discrete logarithm (DLP).

can be derandomized to P/poly (i.e., BPP/poly = P/poly[38]). Note however that BPP/qgenpoly constitutes a physically relevant class, since it only contains problems that can be solved efficiently by classical and quantum computers, as opposed to P/poly, which contains undecidable problems, such as a version of the halting problem. We refer to Supplementary Section 1 for formal definitions of these complexity classes, and an in-depth discussion.

On the other hand, it is easy to show that quantum models (more precisely quantum linear models) can also represent any function in BQP, i.e., all functions that are efficiently computable on a quantum computer. For this, one simply takes a simple encoding of an $n$-bit input $x$:

$$\rho(\boldsymbol{x}) = \bigotimes_{i=1}^{n} X_i^{x_i} |0\rangle \langle 0| X_i^{x_i} \qquad (14)$$

along with an observable

$$O_n = U_n^\dagger Z_1 U_n \qquad (15)$$

specified by an arbitrary $n$-qubit circuit $U_n$ in BQP and the Pauli-$Z$ operator applied on it first qubit. The resulting model $f_n(\boldsymbol{x}) = \mathrm{Tr}[\rho(\boldsymbol{x})O_n]$ can then be used to decide any language in BQP.

Combining these two observations, we get that the proposition "all quantum models are shadowfiable" would imply that BQP ⊆ BPP/qgenpoly ⊆ P/poly, which violates the widely believed conjecture[39] that BQP ⊄ P/poly (see Supplementary Section 4 for a formal proof). To give an example of candidates of non-shadowfiable quantum models, the discrete logarithm $\log_g x \bmod p$ (or even one bit of it) is provably in BQP but is not believed to be in P/poly. Therefore, a model that could be used to compute the discrete logarithm (e.g., the quantum model of Liu et al.[16]) is likely not shadowfiable.

## Discussion

In this work, we examined the class of quantumly trainable, classically evaluatable models we refer to as shadow models. Our analysis has shown that these models can be universally captured by a restricted family of quantum linear models, wherein data-encoding and variational operations are flipped compared to conventional quantum models. Furthermore, we demonstrated that shadow models belong to an intriguing complexity class, coined BPP/qgenpoly, exhibiting superiority over classical models (in BPP) but inferiority to fully quantum models (in BQP), based on prevalent complexity theory assumptions.

By presenting shadows models as flipped linear models, we illustrated how shadow tomography protocols could be applied

straightforwardly to construct shadow models in practice. Yet, it is important to note a crucial distinction between a shadow tomography scenario and a shadow model: in the latter, one has control over the quantum state intended for shadowing. This distinction introduces new possibilities for devising 'state-aware' shadow tomography protocols aimed at constructing shadow models. This could potentially alleviate some of the limitations of current classical shadow protocols.

Considering our findings on learning separations, we identified a noteworthy characteristic of shadow models: their ability to quantumly compute useful advice for a classical evaluation algorithm, enabling them to tackle otherwise classically intractable tasks. The example we presented, based on trap-door functions, readily allows for such constructions, but it remains somewhat contrived. Exploring similar constructions for physically relevant problems, such as predicting ground state properties of complex quantum systems, would be an intriguing avenue for future research.

## Data availability
Data sharing is not applicable to this paper as no datasets were generated or analyzed during the current study.

## References

1.  Biamonte, J. et al. Quantum machine learning. *Nature* **549**, 195 (2017).
2.  Dunjko, V. & Briegel, H. J. Machine learning & artificial intelligence in the quantum domain: a review of recent progress. *Rep. Prog. Phys.* **81**, 074001 (2018).
3.  Schuld, M. & Petruccione, F. *Supervised Learning With Quantum Computers* 1st edn, Vol. 287 (Springer, 2018).
4.  Benedetti, M., Lloyd, E., Sack, S. & Fiorentini, M. Parameterized quantum circuits as machine learning models. *Quant. Sci. Technol.* **4**, 043001 (2019).
5.  Cerezo, M. et al. Variational quantum algorithms. *Nat. Rev. Phys.* **3**, 625 (2021).
6.  Bharti, K. et al. Noisy intermediate-scale quantum algorithms. *Rev. Mod. Phys.* **94**, 015004 (2022).
7.  Schuld, M. & Killoran, N. Quantum machine learning in feature hilbert spaces. *Physical review letters* **122**, 040504 (2019).
8.  Schuld, M., Bocharov, A., Svore, K. M. & Wiebe, N. Circuit-centric quantum classifiers. *Phys. Rev. A* **101**, 032308 (2020).
9.  Liu, J.-G. & Wang, L. Differentiable learning of quantum circuit born machines. *Phys. Rev. A* **98**, 062324 (2018).
10. Jerbi, S., Gyurik, C., Marshall, S., Briegel, H. & Dunjko, V. Parametrized quantum policies for reinforcement learning. *Adv. Neural Inf. Processing Syst.* https://doi.org/10.48550/arXiv.2103.05577 (2021).
11. Skolik, A., Jerbi, S. & Dunjko, V. Quantum agents in the gym: a variational quantum algorithm for deep q-learning. *Quantum* **6**, 720 (2022).
12. Havlíček, V. et al. Supervised learning with quantum-enhanced feature spaces. *Nature* **567**, 209 (2019).
13. Zhu, D. et al. Training of quantum circuits on a hybrid quantum computer. *Sci. Adv.* **5**, eaaw9918 (2019).
14. Peters, E. et al. Machine learning of high dimensional data on a noisy quantum processor. *npj Quant. Inform.* **7**, 161 (2021).
15. Haug, T., Self, C. N. & Kim, M. Quantum machine learning of large datasets using randomized measurements. *Mach. Learn Sci.-Technol.* **4**, 015005 (2023).
16. Liu, Y., Arunachalam, S. & Temme, K. A rigorous and robust quantum speed-up in supervised machine learning. *Nat. Phys.* **17**, 1013–1017 https://doi.org/10.1038/s41567-021-01287-z (2021).
17. Gyurik, C. & Dunjko, V. On establishing learning separations between classical and quantum machine learning with classical data. *arXiv* https://arxiv.org/abs/2208.06339 (2022).

18. Gyurik, C. & Dunjko, V. Exponential separations between classical and quantum learners. *arXiv* https://arxiv.org/abs/2306.16028 (2023).

19. Schreiber, F. J., Eisert, J. & Meyer, J. J. Classical surrogates for quantum learning models. *Phys. Rev. Lett.* **131**, 100803 (2023).

20. Landman, J., Thabet, S., Dalyac, C., Mhiri, H. & Kashefi, E. Classically approximating variational quantum machine learning with random fourier features. *arXiv* https://arXiv:2210.13200 (2022).

21. Schuld, M., Sweke, R. & Meyer, J. J. Effect of data encoding on the expressive power of variational quantum-machine-learning models. *Phys. Rev. A* **103**, 032430 (2021).

22. Jerbi, S. et al. Quantum machine learning beyond kernel methods. *Nat. Commun.* **14**, 517 (2023).

23. Huang, H.-Y., Kueng, R. & Preskill, J. Predicting many properties of a quantum system from very few measurements. *Nat. Phys.* **16**, 1050 (2020).

24. Bertoni, C. et al. Shallow shadows: Expectation estimation using low-depth random clifford circuits. *Phys. Rev. Lett.* https://journals.aps.org/prl/accepted/f9079Y0cPf41d490530c11a3d770886929e5b63ec (2024).

25. Wan, K., Huggins, W. J., Lee, J. & Babbush, R. Matchgate shadows for fermionic quantum simulation. *Commun. Math. Phys.* **404**, 629–700 https://doi.org/10.1007/s00220-023-04844-0 (2023).

26. Hu, H.-Y., Choi, S. & You, Y.-Z. Classical shadow tomography with locally scrambled quantum dynamics. *Phys. Rev. Res.* **5**, 023027 (2023).

27. Schuld, M. Supervised quantum machine learning models are kernel methods. *arXiv* https://arxiv.org/abs/2101.11020 (2021).

28. Cramer, M. et al. Efficient quantum state tomography. *Nat. Commun.* **1**, 149 (2010).

29. Caro, M. C., Gil-Fuster, E., Meyer, J. J., Eisert, J. & Sweke, R. Encoding-dependent generalization bounds for parametrized quantum circuits. *Quantum* **5**, 582 (2021).

30. Gyurik, C. et al. Structural risk minimization for quantum linear classifiers. *Quantum* **7**, 893 (2023).

31. Servedio, R. A. & Gortler, S. J. Equivalences and separations between quantum and classical learnability. *SIAM J. Comput.* **33**, 1067 (2004).

32. Sweke, R., Seifert, J.-P., Hangleiter, D. & Eisert, J. On the quantum versus classical learnability of discrete distributions. *Quantum* **5**, 417 (2021).

33. Kearns, M. J. & Vazirani, U. *An Introduction to Computational Learning Theory*, 222 (MIT Press, 1994).

34. Shor, P. W. Polynomial-time algorithms for prime factorization and discrete logarithms on a quantum computer. *SIAM Rev.* **41**, 303 (1999).

35. Alexi, W., Chor, B., Goldreich, O. & Schnorr, C. P. Rsa and rabin functions: certain parts are as hard as the whole. *SIAM J. Comput.* **17**, 194 (1988).

36. Casas, B. & Cervera-Lierta, A. Multidimensional fourier series with quantum circuits. *Phys. Rev. A* **107**, 062612 (2023).

37. Grover, L. K. A fast quantum mechanical algorithm for database search. In *Proc. 28th Annual ACM Symposium on Theory of Computing* 212–219 (ACM, 1996).

38. Adleman, L. Two theorems on random polynomial time. In *19th Annual Symposium on Foundations of Computer Science (SFCS 1978)* 75–83 https://doi.org/10.1109/SFCS.1978.37 (IEEE, 1978).

39. Aaronson, S. & Arkhipov, A. The computational complexity of linear optics. In *Proc. 43rd Annual ACM Symposium on Theory of Computing* 333–342 (ACM, 2011).

## Author contributions
The project was conceived by S.J., S.C.M., and V.D. The theoretical aspects of this work were developed by S.J., C.G., R.M., and V.D. All authors (S.J., C.G., S.C.M., R.M., and V.D.) contributed to technical discussions and writing of the manuscript.

## Competing interests
The authors declare no competing interests.

## Additional information
**Supplementary information** The online version contains supplementary material available at https://doi.org/10.1038/s41467-024-49877-8.

**Correspondence** and requests for materials should be addressed to Sofiene Jerbi.

**Peer review information** *Nature Communications* thanks the anonymous reviewer(s) for their contribution to the peer review of this work. A peer review file is available.

**Reprints and permissions information** is available at http://www.nature.com/reprints

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.