

Opportunities and limitations of explaining quantum machine learning

Elies Gil-Fuster,^{1,2} Jonas R. Naujoks,² Grégoire Montavon,^{3,4} Thomas Wiegand,^{2,5,4} Wojciech Samek,^{2,5,4} and Jens Eisert^{1,2,6}

¹Dahlem Center for Complex Quantum Systems, Freie Universität Berlin, 14195 Berlin, Germany

²Fraunhofer Heinrich Hertz Institute, 10587 Berlin, Germany

³Department of Mathematics and Computer Science, Freie Universität Berlin, 14195 Berlin, Germany

⁴BIFOLD – Berlin Institute for the Foundations of Learning and Data, 14109 Berlin, Germany

⁵Department of Electrical Engineering and Computer Science, Technische Universität Berlin, 14109 Berlin, Germany

⁶Helmholtz-Zentrum Berlin für Materialien und Energie, 14109 Berlin, Germany

A common trait of many machine learning models is that it is often difficult to understand and explain what caused the model to produce the given output. While the explainability of neural networks has been an active field of research in the last years, comparably little is known for quantum machine learning models. Despite a few recent works analyzing some specific aspects of explainability, as of now there is no clear big picture perspective as to what can be expected from quantum learning models in terms of explainability. In this work, we address this issue by identifying promising research avenues in this direction and lining out the expected future results. We additionally propose two explanation methods designed specifically for quantum machine learning models, as first of their kind to the best of our knowledge. Next to our pre-view of the field, we compare both existing and novel methods to explain the predictions of quantum learning models. By studying explainability in quantum machine learning, we can contribute to the sustainable development of the field, preventing trust issues in the future.

I. INTRODUCTION

Concomitant with the advent of enormously large databases, the striking development of available computing power, and a rapid increase in the pace of research on data-driven approaches, systems of *artificial intelligence* (AI) and *machine learning* (ML) are reaching levels of performance unseen before, often outperforming human capabilities [1, 2]. From chat systems passing the Turing test [3], to achieving super-human performance in strategy games [4], and becoming ubiquitous in speech [5], image [6], and other pattern recognition tasks, AI algorithms have changed our world.

Given the importance of AI algorithms to our modern lives, researchers have been thinking intensely recently whether one could possibly improve such algorithms by resorting to *quantum computers*. While presently-available quantum devices are still noisy and medium-sized, the rate of development of quantum technologies prompts us to imagine a near future with large-scale quantum computers. On some paradigmatic—albeit not yet practical—mathematically well-defined problems, quantum computers come within reach or outperform the fastest classical supercomputers available to date [7, 8]. Generically, a key question is still whether quantum computers can achieve better performance in learning tasks than classical algorithms, and how this would manifest in practice. This is the core question of the aspiring field of *quantum machine learning* (QML) [9–12]. Such quantum advantages—as this state of affairs is often referred to—could refer to advantages in sample complexity, in computational complexity, or in generalization. For highly structured problems, super-polynomial quantum advantages have actually already been proven [13–15], and it can only be a matter of time to have further evidence available to what extent quantum algorithms might possibly help dealing with real-world data.

One disadvantage of state-of-the-art ML models is the difficulty to understand the model’s behavior as an emergent phenomenon from inspecting the individual constituent compo-

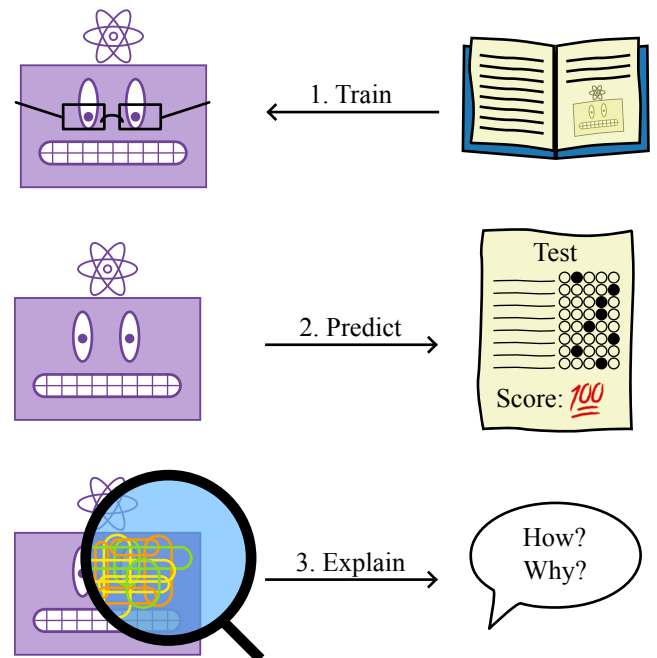


Figure 1. Framework of eXplainability in *quantum machine learning* (XQML) as presented in this work.

nents [16, 17]. This hinders broader ML adoption in applications where human lives and social justice are at stake. It is the goal of *explainable AI* (XAI) [18] and, more broadly, of reliable AI, to mitigate the potential negative impact of ML in society by unveiling the causes behind the behavior of modern learning models. Notions of XAI provide answers to the question of what it actually means to explain a learning model. One here asks questions of the type, *what is the role of each of the parameters in the model?*, *what pixels in this image are responsible for it to be labelled as a cat?*, or *how can we design learning models which remain human-interpretable through-*

out learning? It is not a single method, but rather a portfolio of tools that capture various aspects of this question.

While this field of research is for good reason becoming well-established for classical AI algorithms [18–20], the same cannot be said for quantum algorithms. In fact, few steps have been taken so far to bring transparency to quantum ML [21–24]. This lack of a machinery gives rise to a grave omission, however: All limitations of black-box classical ML (e.g., in terms of trustworthiness) also exist in black-box QML. Hence, it is important to find out what feature lets a quantum algorithm make a certain prediction, to “open the black box”, metaphorically speaking.

There is one major difference between the classical and quantum settings in terms of explainability. Namely, QML research still has not realized the main goal in the field: to solve practically relevant problems faster or otherwise better than its classical counterpart. This represents a timely opportunity for developing *explainable quantum machine learning* (XQML) techniques: we have the chance to design QML models that are inherently explainable from the get-go. Conversely, the development of high-performing models in classical ML predates the development of the field of XAI, resulting in tensions between accuracy and interpretability [25]. The cartoon in Fig. 1 introduces QML as a three-step process: after training and observing good tests results, we need to explain the behavior of the quantum learning agent.

In this work, we aim at presenting first steps towards establishing such a framework. Our main contribution is then two-fold:

- We offer a broad perspective of the incumbent field of explainability for QML: We do so by introducing and reviewing methods of explainability and what is specifically different for QML.
- We also propose two concrete explanation techniques designed specifically for present-day QML models, based on *parametrized quantum circuits* (PQCs).

II. PRELIMINARIES

In this section, we attempt to bridge the fields of (classical) explainable AI and QML. We first lay out the formalism of parametrized quantum circuits, which are the workhorse in most of present-day QML. Our introduction focuses on the essential differences between classical and quantum ML that we inherit from the principles of quantum mechanics. This way, we do not dwell extensively on the QML practitioner’s point of view. Next we present explainability holistically, singling out the types of questions one may ask in order to explain the behavior of a learning model. Here we start very broad and become progressively narrow until we reach the so-called post-hoc local explanations, which are the area of explainability we concentrate later on, in Section III.

A. Quantum function families for machine learning

In this work we restrict ourselves to learning tasks with classical data. The predominant approach to QML for classical data relies on a view of quantum computing that is reminiscent of probabilistic classical computers. A common reading of quantum mechanics is as a generalized theory of probability, where qubits play the role of probabilistic bits. We keep the discussion here at a high level, so that an understanding of physics is not required to follow the exposition.

A quantum state composed of many qubits can be thought of as a generalized probability distribution over bit-strings. Probability distributions over n -bit-strings can be represented as diagonal matrices

$$D_p = \sum_{i \in \{0,1\}^n} p_i |i\rangle\langle i|, \quad (1)$$

where the notation $|i\rangle\langle j|$ is the Euclidean basis of the vector space of matrices. Given that its diagonal is real, it follows that the matrix is Hermitian $D_p = D_p^\dagger$. Moreover it is also *positive semi-definite* (PSD) $D_p \geq 0$, since all entries are non-negative. Finally, due to normalization of any probability distribution p , it has unit trace $\text{tr}(D_p) = 1$. The step from classical probability distributions to *quantum states* is to just remove the requirement that the matrix be diagonal. This way, a quantum state ρ over n qubits is represented by a $2^n \times 2^n$ -dimensional matrix (same as D_p), which is Hermitian, PSD, and has unit trace. The off-diagonal entries of the matrices are in general complex valued, and we refer to them as *coherences*. Classical probability distributions correspond to the special cases of quantum states with diagonal matrices. Then, a quantum state stores statistical information on its diagonal. We say it stores statistical information because the diagonal always represents a probability distribution, and we can sample this distribution by physically performing a measurement in the computational basis. The information in the off-diagonal corresponds to other correlations among different bit-strings that would influence further transformations of the state.

In simple terms, a quantum circuit is a sequence of operations each of which takes a quantum state as input and returns another quantum state as output. Transformations of quantum states can be thought of as analogues of transformations of probability distributions. For example, a doubly stochastic matrix S is a linear transformation that maps probability distributions onto other probability distributions $p' = Sp$. Similarly, the linear transformations that map quantum states onto quantum states are called *quantum channels*. In general, these are completely positive, trace preserving maps, and, in particular, the action of any stochastic matrix can be seen as a quantum channel acting on the corresponding matrix $D_p \mapsto D_{p'}$ as a convex combination of permutations. In general quantum channels can correspond to exotic transformations of states, but in practice we mainly consider those that we can hope to implement on quantum hardware. Therefore, throughout this work, we only consider unitary transformations, so quantum channels that act as $\rho \mapsto U\rho U^\dagger$, where U is a unitary matrix, fulfilling $U^\dagger U = \mathbb{I}$ by definition.

A common parametrization of the set of all unitary matrices is via a collection of rotation angles $\vartheta = (\vartheta_1, \dots, \vartheta_M) \in [0, 2\pi)^M$. Then, we usually talk about parametrized transformations $\rho \mapsto \rho'(\vartheta) = U(\vartheta)\rho U^\dagger(\vartheta)$. Unitary transformations are also referred to as *quantum gates* in the context of quantum computing, as they are complex-valued generalizations of invertible logic gates.

Finally, a real-valued function that takes bit-strings $b: \{0, 1\}^n \rightarrow \mathbb{R}$ as input provides a natural way to extract information from a probability distribution via its expectation value $\langle b \rangle_p = \sum_{i \in \{0, 1\}^n} b(i) p_i \in \mathbb{R}$. We call

$$B = \sum_{i \in \{0, 1\}^n} b(i) |i\rangle\langle i|, \quad (2)$$

which is again a real-valued, diagonal, Hermitian matrix (just not necessarily unit trace nor PSD). Using this representation, the expectation value adopts the form of the Hilbert-Schmidt norm of the matrices $\langle b \rangle_p = \text{tr}\{D_p B\}$. The way via which we retrieve information from a quantum state ρ is fully analogous: Given an observable \mathcal{M} (a complex-valued Hermitian matrix), we can compute its expectation value relative to the state as $\langle \mathcal{M} \rangle_\rho = \text{tr}\{\rho \mathcal{M}\}$, which is always real-valued due to both ρ and \mathcal{M} being Hermitian.

A typical quantum computation then takes the form of a three-step process:

1. Start from a fixed (easy-to-prepare) initial state ρ_0 .
2. Apply a *quantum circuit* specified as a sequence of (parametrized) unitary gates $U_1(\vartheta_1), \dots, U_L(\vartheta_L)$, each of which acting as $\rho_j = U_j(\vartheta_j)\rho_{j-1}U_j^\dagger(\vartheta_j)$, and resulting in the final state $\rho_L(\vartheta_1, \dots, \vartheta_L)$.
3. Estimate the expectation value of a fixed (easy-to-measure) observable¹ \mathcal{M}_0 , obtaining the real value $\text{tr}\{\rho_L(\vartheta_1, \dots, \vartheta_L)\mathcal{M}_0\}$.

The real-valued result of this quantum computation is

$$f(\rho_0, (U_j(\vartheta_j))_j, \mathcal{M}_0) = \text{tr}\{\rho_L(\vartheta_1, \dots, \vartheta_L)\mathcal{M}_0\}. \quad (3)$$

We call a *parametrized quantum circuit* (PQC) a fixed choice of initial state ρ_0 , parametrized gates $(U_j)_j$ (not their parameters), and final observable \mathcal{M}_0 . Each PQC gives rise to a real function

$$(\vartheta_j)_j \mapsto g_{\rho_0, (U_j)_j, \mathcal{M}_0}((\vartheta_j)_j) = f(\rho_0, (U_j(\vartheta_j))_j, \mathcal{M}_0) \quad (4)$$

of the parameters. We abuse notation slightly and still use f instead of $g_{\rho_0, (U_j)_j, \mathcal{M}_0}$ when the specific PQC is irrelevant or clear by context. We refer to functions arising from PQCs in

¹ As a technical aside, talking about an evolved state with a fixed observable is called the *Schrödinger picture*. We could have analogously talked about a fixed state and an evolved observable, placing ourselves in the *Heisenberg picture*. Considering the space in-between, where both the state and the observable are evolved via some gates each is sometimes called the *interaction picture*, and we exploit it later.

this way as *quantum functions*. Appendix A contains a brief description of an analogous family of classical functions defined from parametrized probability distributions.

To perform machine learning it is not enough to have a single function, though. If we call \mathcal{X} the data domain, and \mathcal{Y} the label co-domain, we must furnish a family of functions $\mathcal{F} \subseteq \mathcal{Y}^{\mathcal{X}}$, to be used as the hypothesis class. One straightforward approach would be to, given a PQC, consider a set of possible initial states $\{\rho(x) | x \in \mathcal{X}\}$ instead of a single one, or equivalently a data-dependent observable $\mathcal{M}(x)$. This has been a popular approach since the start of PQC-based QML, and it knocks on the door of the fundamental question “what is the best way to upload classical data into a quantum computation?”. In practice, having a data-dependent state amounts to keeping the structure introduced above, and just fixing the rotation angles of the first gates to be equal to the input data. One commonly distinguishes between an *encoding* gate $E(x)$, and the *trainable* (or *variational*) gates $U_j(\vartheta_j)_j$. With this notation the data-dependent state is $\rho(x) = E(x)\rho_0 E^\dagger(x)$. We refer to $V_L(\vartheta) = U_L(\vartheta_L) \dots U_1(\vartheta_1)$ as the quantum gate containing all the trainable gates, and we call Θ the set of all possible parameter values $\vartheta \in \Theta$. In this prescription, a fixed PQC (where the encoding gate is also fixed as of step 2 above) gives rise to a function

$$f(x, \vartheta) = \text{tr}\{V(\vartheta)\rho(x)V^\dagger(\vartheta)\mathcal{M}_0\} \quad (5)$$

of both the input x and the variational parameters ϑ . The division between encoding and trainable gates allows us to consider the set of all functions

$$\mathcal{F} := \{x \mapsto f(x, \vartheta) | \vartheta \in \Theta\} \quad (6)$$

that arise from different parameter values for a given PQC. In general terms, these are the quantum function families that we consider for QML. Each choice of PQC amounts to a different hypothesis class in the exact same way as different neural network (NN) architectures result in different ML models. Specifically, the milestone works [26, 27] showed that PQCs give rise to generalized trigonometric polynomials. That is, the parametrized function families realized by such circuits are always of the form

$$f_\vartheta(x) = \sum_{\omega \in \Omega} (a_\omega(\vartheta) \cos\langle \omega, x \rangle + b_\omega(\vartheta) \sin\langle \omega, x \rangle). \quad (7)$$

Here, the values of the coefficients $(a_\omega(\vartheta), b_\omega(\vartheta))_\omega$ depend only on the trainable parameters ϑ of the PQC, in a computationally intricate way. Further, the frequency spectrum Ω depends only on which are the encoding gates. Here $\langle \cdot, \cdot \rangle$ denotes the usual Euclidean inner product of vectors.

The study of PQC-based QML models, which we consider in this work, has been pursued within the quest for practical quantum advantages in meaningful learning tasks, as potential quantum analogues to neural networks. Although both approaches rely on parametrized building blocks and share a layered structure, the principles of quantum mechanics prevent much deeper similarity:

- The *no-cloning-theorem* – that shows that unknown quantum states cannot be copied so that both copies

are statistically indistinguishable from the input – prevents parametrized gates from fanning out the same way as artificial neurons do. Namely, the signal a classical neuron transmits can be copied freely and used as input for multiple neurons in the next layer. For PQCs, a leading framework to replicate this feature is data re-uploading [28]².

- In general, the matrix representation of quantum states requires an exponential dimension, which means we cannot efficiently store the intermediate state of the computer at each step.
- Storing intermediate information (even if it is not the entire state) necessarily destroys the computation, which prevents information re-use and the existence of fast training algorithms like back-propagation [29].
- In practice, we estimate expectation values via samples from the quantum computer. This means that our read-out precision is only polynomial in time, not exponential, which results in distinguishability problems and can mask the signal inside the statistical noise. This problem is discussed also in the context of concentration phenomena and vanishing gradients [30, 31].
- The application of unitary gates is a linear transformation (non-linear transformations are forbidden in quantum mechanics). The only possible source of non-linearity comes from the parametrization of the gates $\alpha \mapsto U(\alpha)$, and not from the nested composition of information-processing layers. This allows us to understand the expressivity of PQC-based QML models as trigonometric series [27].
- The connectivity graph of the underlying quantum computing platform can severely restrict the unitary gates we can implement in practice. In classical NNs sparsely-connected layers are used among other reasons because doing so has been empirically found to boost performance. Conversely, for PQC designs to remain realistic we should actually only consider unitary gates that act on at most 2 qubits at once, and ideally they should be adjacent. In spite of the field of quantum compilation being by now mature, the hardware limitations in the short- and mid-term must be kept in mind.

There exist ways to circumvent some of these “problems” by considering only restricted families of quantum circuits. Yet, these characteristics are fundamental to quantum mechanics: by aiming to remove them even partially, we may inadvertently fall in the regime of *efficient classical simulability*, in which there can be no quantum advantage [32–40]. Indeed, we must take into account the classical and quantum time and computational complexity of both the QML models themselves and the explainability methods we consider later.

B. Explainability

In addition to building models that perform accurate predictions, it is often desirable to ensure they operate with some level of transparency, in particular, ensuring that their decisions can be explained in a way that is understandable for a human. This aspect has been treated extensively within the field of *explainable AI* (XAI) [18, 20, 41]. A central example motivating the usefulness of explainability is detecting when a learning model is performing well “for the wrong reasons”, the so-called clever Hans effect [42, 43] or “shortcut learning” [44]. A flawed prediction strategy, e.g., one relying on spuriously correlated artifacts, may stop working on new data, exposing the environment in which the model operates to great risks (see, e.g., Refs. [43, 45]). Explainable AI techniques can precisely pinpoint these flawed decision strategies, and also provides a starting point to produce more robust models by aligning the model with human expert knowledge [46–48]. The same Explainable AI methods can also enable a reverse alignment where a ML model trained on vast amounts of data can serve as an exploration tool for humans, enabling for example, the generation of new insights or scientific hypotheses [49, 50].

The field of XAI encompasses a multitude of diverse techniques addressing different levels of interpretability as well as models and data to which they apply. Explanation techniques can be organized along several conceptual axes:

a. Global versus local explanations. A local explanation focuses on unraveling a model’s prediction proper to a single data point. Such explanations then result in properties defined for the input pattern, like the relevance of each of its dimensions. A global explanation on the other hand is one that aims at elucidating a given model’s overarching reasoning strategy. This includes recognizing prototypical behavior of the model without inherently providing an interface to understand a given single data point. A representative global explanation method is *spectral relevance analysis* [42, 51], which looks at the behavior of a model on a dataset by aggregating the output of a large collection of local explanations, e.g., into clusters that summarize the complex multifaceted decision strategy of the model.

b. Feature-wise versus concept-wise explanations. Another way of distinguishing between explanation techniques is with respect to what are the actual units of explanation. While classical explanations typically explain decision of the model in terms of their input features (e.g., pixels [52]) or combinations of them (see, e.g., Ref. [53]), a more recent line of research has sought to support the explanation by more abstract features, so-called concepts, which can be, for example, activations in intermediate layers, or linear transformations of those activations [48, 54–56]. Concept-based explanations, because of their intrinsic lower dimensionality, are also more amenable to a global (dataset-wide) analysis of the model.

c. Ante-hoc versus post-hoc explanations. An additional line of disambiguation of explainability techniques concerns the *stage* at which explanations are generated. The field of ante-hoc explainability deals with the design and study of learning models that are intrinsically interpretable, specif-

² In data re-uploading, the input is uploaded several times to introduce a redundancy which is present by default in NNs.

ically, they embed units of activation that are directly inspectable and that carry intrinsic meaning. This encompasses generalized additive models [57], decision trees, or neural networks with specific topologies (e.g., global average pooling layers [58] or attention layers [59]). On a related note, data-driven physical theories, e.g., implemented via neural ordinary differential equations [60], constitute another instance of ante-hoc explanations: For example, before we fit the viscosity parameter to the measured data, we already know that the role of that parameter will be the viscosity coefficient, and thus the viscosity model is inherently interpretable. Enforcing such a structure in a learning model limits its capabilities compared to black-box models. This concept is known as the *performance-explainability trade-off* [20]. Nowadays, the most powerful ML models, especially deep neural networks, fall under the category of opaque or black-box models: they are powerful but their inner workings are not easily interpretable. In this case, post-hoc explainability is used. A post-hoc explanation is extracted from an already trained model.

d. Attribution versus synthesis. We refer to an explanation as *attribution* when the focus is on which features are measurably responsible for the predicted output. Many attribution methods (see, e.g., Refs. [52, 61, 62]) are designed to assign a score to a particular feature that can be interpreted as the share by which that feature contributes to the output of the model. In contrast, *synthesis* methods rely on generating new data points which carry insight, either on their own or in relation to the original data point. For example, given a correctly classified input, one can generate a counterfactual [63], specifically, the smallest natural perturbation to that input that causes a change in the predicted class. One could also, without relying on any particular data point, synthesize a so-called class prototype that maximizes the score for a given class [64]. This synthesis approach to explanation is also the basis of the *DeepDream* technique [65] (which was adapted to quantum circuits in Ref. [22]).

C. Post-hoc local attribution methods

We focus specifically on *post-hoc local attribution* methods, in which explanations consist of assigning a relevance score to each input dimension of a given pattern. The use cases we present later all belong to this class of explanations. This presupposes that the input features or relevance attributions thereof are themselves interpretable to the end user. In a classification setting, an ML model typically outputs a score for each of the K classes $h(x) = (h_1(x), \dots, h_K(x))$. To explain a single input, though, we call f the score for the class that we would like to investigate. Often, one chooses the class that corresponds to the true label for each input. So, explicitly, the human-interpretable property we explain in this case is “how much does each input dimension contribute to the correct class?” For local attribution methods, the final explanation is sometimes called a *relevance heat-map* (as for image-recognition tasks that is what it looks like). Many attribution techniques have been proposed, which we organize in

the following into three types of approaches.

a. Occlusion-based approaches. A first approach to attribution consists of scoring input features according to the effect of adding or removing them on the output of the model. This approach is instantiated by the method of *Shapley values* (SV) [66], which originates from coalition game theory. In its original formulation, the method addresses the question of assigning the individual contribution of a given set of players to a global payoff function. The method comes with strong theoretical foundation (it has a unique axiomatic characterization [66]) and has been extended to attribute a machine learning model (see, e.g., Refs. [62, 67]), where the players forming the coalition are the input features, and the payoff to be distributed correspond to the output of the model to be explained.

In this work, we consider the *baseline Shapley values* approach [68], where such features are set to a baseline \tilde{x} . Starting from the original datum x , and given the baseline \tilde{x} , one systematically replaces each possible combination of input components of x by those of \tilde{x} . Such an explanation is formalized as

$$E_i(x) = \sum_{S \subseteq ([d] \setminus \{i\})} \frac{|S|(d-1-|S|)!}{d!} (f(x_{S \cup \{i\}}) - f(x_S)), \quad (8)$$

where $[d] = \{1, \dots, d\}$, and for any set $A \subseteq [d]$, we use x_A to denote the modification of x whose component indices in A are left untouched. So the i th component of x_A is x_i if $i \in A$ and \tilde{x}_i otherwise. Computing Shapley values exactly requires computational resources that scale combinatorially, so in practice they are most often replaced by a sampling approximation (like *Shapley value sampling* [62]). There exist variations for specific model classes where the computation of Shapley values scales polynomially [69, 70].

b. Gradient-based approaches. An alternative to evaluating the effect of removing each features individually consists of relying on the gradient of the model, which can be interpreted as a local effect of each individual feature. Gradients can usually be efficiently computed, for instance via the backpropagation algorithm for neural networks [71]. A particular way of turning a gradient into an attribution is to use it to perform a first-order Taylor expansion of the function f around some baseline \tilde{x} , evaluate it at the point of interest x , and identifying the first-order terms bound to each input dimension, i.e.,

$$E_i(x) = \partial_i f(\tilde{x})(x_i - \tilde{x}_i). \quad (9)$$

For the method to achieve a meaningful attribution of the model output, the first-order terms must dominate over the higher-order terms locally, which essentially requires the function to be locally linear. This often requires to adapt the root point to the actual data point, i.e. define a function $\tilde{x}(x)$. We refer to this explanation method as `Taylor-1`.

Another gradient-based approach to compute explanations is *integrated gradient* (IG) [61], where rather than evaluating the gradient only once at \tilde{x} , we integrate it over a path connecting \tilde{x} and x . If choosing the path to be a straight line, one gets the integral $\int_{[0,1]} \partial_i f(tx + (1-t)\tilde{x})dt$, which we can

substitute to the gradient in the equation above. In practice, the integral is approximated by a discretization. The main advantage of the Integrated Gradients approach over Taylor expansion is its ability to produce meaningful explanation even in the presence of strong non-linearities. This broader applicability of Integrated Gradients comes however at the cost of a higher computational cost.

Another simple yet commonly used gradient-based explanation technique is `grad×input`: $E_i(x) = \partial_i f(x)x_i$. For specific functions, like homogeneous linear functions, `grad×input` can be shown to coincide with the other gradient-based approaches.

c. Propagation-based approaches. This last family of methods differ from the ones above by viewing the prediction not as a function but as the output of a computational graph. It then seeks to attribute the prediction to the input features by reverse propagating in this graph. Such an approach leads to explanations that are specific to the computational graph. This is in contrast to explanations that only require black-box access to the model, which we refer to as *model-agnostic* explanations. The approach is exemplified by the *layer-wise relevance propagation* (LRP) [52] approach. Originally proposed in the context of convolutional neural networks and bag-of-words models, LRP was later extended to a broad range of models (e.g., unsupervised models [72] or transformers [73]). In its basic form, LRP assumes that the function $f(x)$ of interest can be expressed as a sequence of mappings, e.g. in a deep neural network the mapping onto its successive layers:

$$\{a_l \mapsto a_{l+1}\}_{l=1}^{L-1}, \quad (10)$$

where a_1 and a_L are the input and output of the model to explain. Explanation proceeds by iteratively attributing the model output to the diverse layers of the network, starting with the top layers and moving backward towards the input. This layer-wise attribution is implemented by purposely defined propagation rules. Denoting by R_l and R_{l+1} the attribution on the representation a_l and a_{l+1} , respectively, the LRP propagation rule can be generically represented by the collection of maps:

$$\{(a_l, R_{l+1}) \mapsto R_l\}_{l=1}^{L-1} \quad (11)$$

at each layer, defining the redistribution strategy. The process starts with setting $R_L = f(x)$, applying the maps above at each layer in reverse order, and retrieving the vector R_1 , which forms the desired explanation $E(x)$. We illustrate this process with pictorial diagrams in Appendix B. The main advantage of the LRP approach over occlusion and gradient-based methods discussed above is its higher flexibility, allowing to finely control the forming of the explanation at each layer, addressing its specific nonlinearities, so that a robust explanation can be extracted in only a single forward-backward pass. For a more detailed introduction of the LRP technique, including explicit propagation rules, we refer the reader to Ref. [74].

D. Evaluating explanation techniques

Having introduced various approaches to explainability, as well as a variety of technique to achieve post-hoc local attributions, one question remains, namely which explanation technique is most appropriate in a given context. The first and most important criterion to be considered is *explanation correctness*, that is, if the explanation faithfully reflects the underlying strategy the model has used to achieve its prediction. Unlike the prediction function itself, there are typically no ground-truth explanations to be evaluated against. The strategy used by the model to be explained is by definition unknown. Explanation correctness can thus only be assessed via *indirect* means.

A first property an explanation should satisfy is *conservation*: Conservation requires that, if one is to interpret the score assigned to each input features to be the share they contribute to the function output, one should have that $\sum_i E_i(x) \approx f(x)$. Another property an explanation should satisfy is *continuity*, requiring that similar inputs fed into a continuous prediction function should receive similar explanations. A further property to be satisfied is *implementation invariance*, requiring that for two different models implementing the same function f , the explanation for a given data point should be the same or similar. It is noteworthy that the fulfillment of these properties can be assessed to some extent by a direct inspection of the explanation technique without requiring experimentation with a real model trained on real data. For example, one can demonstrate that SV gives rise to continuous explanations, being a weighted sum of continuous functions. Conversely, gradient-based method may produce discontinuous explanation if we consider that the gradient of a continuous prediction function may be discontinuous. An analysis of presented methods in light of these three criteria is given in Table I (left). These formal properties of an explanation may however be insufficient to fully characterize its correctness. For example, a simple explanation technique that would uniformly redistribute the function output to the input features would fill all correctness criteria above, yet fail to highlight features that contribute the most to the prediction.

A common way of assessing the ability of an explanation consists to retrieve the most relevant features from all input features is the “pixel-flipping” evaluation procedure [52, 75]. Pixel-flipping removes features from most to least relevant and keep track of how quickly the output of the model decreases. The faster the decrease the better the explanation. An alternate way of testing a explanation method’s ability to extract correct features is to devise synthetic datasets and models for which the ground-truth explanation is known and test the degree of consistency the explanation technique offer with respect to these ground-truth explanations [76]. We propose specific such metrics in Appendix F, which we use in our numerical experiments in Section IV. Acknowledging the multifaceted nature of the task of evaluating explanations, works like Ref. [77] have contributed series of explanation correctness tests together with their software implementation, which can be used for benchmarking purposes.

Besides explanation correctness, further aspects need to be

Table I. Evaluation of surveyed local attribution methods in terms of formal (data-set unspecific) criteria of explanation correctness and other more general usefulness criteria.

Local Attribution Method	Explanation correctness			
	Conservative?	Continuous?	Implementation invariant?	Computational efficiency?
SV	✓	✓	✓	✗
Taylor-1	✓	✗	✓	✓
IG	✓	✓	✓	✗
grad×input	✗	✗	✓	✓
LRP	✓	✓	✗	✓

considered for an explanation to fulfill its practical purpose. Listing desirable properties for human-interpretable explanations is an enterprise that pre-dates modern machine learning: remarkable work was already laid out in Ref. [78]. Among the list of possible desiderata, one can mention *computational efficiency* which requires that the explanations are computable within ideally one single evaluation of the machine learning model (the computational efficiency of the presented method is listed in Table I. Other desiderata such as human interpretability or sufficiency of the explanation pertain more to the class of explanation techniques, for example, whether the explanation is returned in terms of input features or in terms of more abstract human-readable features, or whether the scoring of those feature (feature attribution) provides sufficient information for delivering useful insights and make the explanation actionable.

E. Explaining PQC-based QML models

We close this section by briefly discussing whether and which explanation techniques can be carried over directly into QML. We start by revisiting the list of fundamental limitations inherited from quantum mechanics from Section II A and, for illustrative purposes, how they relate to the models and properties introduced in Sections II B and II C. We summarize the main points of this section in Table II.

The no-cloning-theorem mostly limits the QML models themselves, and not so much what can be explained about them. Just every time the model is evaluated, a the quantum circuit needs to be ran from scratch, without storing any intermediate information. The exponential dimension of the Hilbert space of quantum states and observables indicates that naive uses of intermediate information may be possible, but at the cost of exponential memory requirements and runtime, which present a problem for scalability. The problem of storing intermediate information exactly on classical memory can be alleviated by considering quantum-time-efficient approximate representations, for which there are several proposals [38, 79–83]. The finite-shot noise is another relevant practical problem: we are only able to evaluate quantum functions (expectation values of PQCs) to precision polynomial in the total runtime, not exponential like in classical computers. The precision we are able to achieve in estimating the model

in turn limits the precision of gradients (this is an important current obstacle for training QML models [30, 31, 84, 85]). Simply put, we are unable to distinguish between exponentially small values and zero. Therefore, all XQML methods that require evaluating gradients might be rendered inapplicable for the common cases where typical gradient magnitudes are exponentially small. Next, the fact that the only source of non-linearity in PQCs is the parametrization of the unitary gates tells us that we shall often be able to exploit the linear nature of quantum computations to our advantage, since explaining linear models is relatively straightforward. All in all, the main factors that play a role in the explainability of QML models are the storage of intermediate information, the finite precision in estimating hypotheses and their gradients, and the fundamentally linear nature of quantum computation.

Further differences arise when comparing specifically QML models to NNs from the point of view of explainability. One key fact is that, even though the graphical depictions of PQCs often resemble (at least structurally) those of NNs, in truth these are objects of essentially different nature. The usual sketches of NNs as subsequent layers of neurons and their connections are actually the computational graph of the model itself. Conversely, the usual sketches of PQCs as subsequent layers of unitary gates applied to several qubits capture the step-by-step physical implementation of these circuits in a laboratory, but not the computational graph of a classical representation of the same process. This difference in computational graphs may already be enough to point out that model-specific techniques that have been derived for NNs shall not be in general directly applicable to PQCs.

An operational difference between QML and NNs results from the different sources of non-linearity. Networks that use *rectified linear units* (ReLU) give rise to functions that are only *piece-wise* analytic (their gradients are only piece-wise continuous), which can give rise to so-called “shattered gradients” for points close to the non-analytic ones. Conversely, QML models are sums of trigonometric functions with bounded frequency, which gives rise to functions that are analytic everywhere. The intuition behind the importance of explanation continuity is that small perturbations of the input should lead to small perturbations of the explanation. The fact that NNs give rise to non-analytic functions results in explanations not being continuous by default. Conversely, all XQML methods considered so far are always continuous because the

Table II. Summary of main differences between NNs and QML models, and their effects on explainability.

	Neural Networks	PQC-based QML models	Effect on explainability
Intermediate information	Stored by default	No cloning theorem and exponential classical memory requirements	Model-specific explanations using intermediate information are more difficult to find for QML models.
Available precision	Exponential	Polynomial	Limited available quantum hypothesis families.
Linearity	Nested composition of non-linear activation functions	Non-linear parametrization of high-dimensional linear operations	The linearity of quantum mechanics can result in straightforward explainability.
Computational graph	Usual structural depiction	Not usual structural depiction	In general not possible to directly port model-specific explanations from NNs to QML. For instance some divide-and-explain ideas like common forms of LRP cannot be immediately applied.
Continuity of gradients	Piece-wise continuous	Analytic everywhere	Continuity of explanations is less problematic for QML than for NNs.

labeling functions themselves are analytic.

In extending explanation methods to quantum models, we can also discuss the expected efficiency considerations. In principle, black-box explanations that only require evaluating the hypotheses should not incur efficiency issues, the total quantum complexity should be the total number of calls to the quantum model times the complexity of evaluating the model. Reasonably speaking, we shall only attempt to produce explanations for quantum circuits which are efficient to evaluate in the first place, so the complexity of producing black-box explanations for quantum models should also be computationally efficient. It should also not be a problem if the black-box model requires evaluating gradients. Issues might arise from the required precision, as discussed above, but that is a different notion of efficiency. For model-specific explanations we further use information about the structure of the model. If the used information is only the succinct description of the circuit, that should also not be problematic. Nevertheless, if the information required to evaluate the model-specific explanation involves storing intermediate information, that might incur exponential classical space overheads. Below we propose two novel model-specific explanations that illustrate precisely this distinction: we provide first a black-box explanation that is quantum efficient, and second an explanation that requires storing intermediate states in classical memory, which uses intractable amounts of space. As discussed, there could exist model-specific XQML methods which are efficient with a smart strategy using the same quantum circuit but with mid-circuit measurements. Such strategies might result in explanations that display a quantum-classical separation, if one can prove that the same explanation could not be realized without access to the quantum circuit.

F. Related work on XQML

So far, a few incursions have already been performed toward XQML. For instance Ref. [21] performed a numerical

study for simple PQCs using explanations like IG and SV and recorded the effects of noise in their explanation performance. Also, Ref. [23] extended the definition of Shapley values to be well-defined for functions defined as expectation values of random variables and then proposed several applications of the resulting SV -like explanations. The authors of this work did not focus only on QML applications, like local attribution methods, but also they used the mindset of coalition games to explain other properties of PQCs, like the effect of each individual trainable gate in the circuit. Finally, Refs. [22, 24] proposed approaches to interpretability (related to explainability). In Ref. [22] the technique of *deep dreaming* is used to improve the design of quantum circuits and to extract human-interpretable properties. Improving on a popular method consisting on building a local interpretable surrogate of the model called LIME [86], Ref. [24] develops a method in which PQCs not only produce a prediction, but also a confidence value for it.

III. NOVEL XQML APPROACHES

Here we propose two novel explanation techniques designed specifically for PQC-based QML models. To the best of our knowledge, these are the first model-specific explanations for quantum learning models. The first of the two, which we call $\text{Taylor-}\infty$ is a black-box explanation that exploits the Fourier picture of PQCs [27]. The second, which we call *quantum layerwise relevance propagation* (QLRP) is our proposal to adopt the *divide and explain* philosophy behind the LRP explanation technique [87]. In this section we introduce the main ideas behind these techniques, a full derivation together with implementation-oriented explanation can be found in Appendices C, D, and E.

A. Taylor- ∞ explanation

To arrive at Taylor- ∞ we recall from Refs. [26, 27] that PQC's where each input component is encoded only once give rise to degree-1 trigonometric polynomials. That is, the functions realized by such circuits are always of the form

$$f_{\vartheta}(x) = \sum_{\omega \in \Omega} (a_{\omega}(\vartheta) \cos\langle \omega, x \rangle + b_{\omega}(\vartheta) \sin\langle \omega, x \rangle), \quad (12)$$

where the values of the coefficients $(a_{\omega}(\vartheta), b_{\omega}(\vartheta))_{\omega}$ depend only on the trainable parameters of the PQC. Further the frequency spectrum $\Omega \subseteq \{0, \pm 1\}^d$ is such that $\Omega \cup -\Omega = \{0, \pm 1\}^d$, $\Omega \cap -\Omega = \{0\}^d$. Here $\langle \cdot, \cdot \rangle$ denotes the usual Euclidean inner product of vectors.

We propose the explanation Taylor- ∞ as a generalization of Taylor-1 in which we take not only the first-order terms $\partial_i f(\tilde{x})(x_i - \tilde{x}_i)$, but rather all single-component contributions $\partial_i^k f(\tilde{x})(x_i - \tilde{x}_i)^k$. Indeed, the Taylor- ∞ explanation is based on the expression

$$f(x) = f(\tilde{x}) + \sum_{i=1}^d \left(\sum_{k=1}^{\infty} \frac{\partial_i^k f(\tilde{x})(x_i - \tilde{x}_i)^k}{k!} \right) + \varepsilon. \quad (13)$$

Now ε contains only higher-order crossed terms, including partial derivatives with respect to different components. We denote the infinite series by $T_i(x, \tilde{x})$, and exploiting usual trigonometric identities, we reach the defining formula

$$T_i(x, \tilde{x}) = \sin(x_i - \tilde{x}_i) \partial_i f(\tilde{x}) + (1 - \cos(x_i - \tilde{x}_i)) \partial_i^2 f(\tilde{x}) \quad (14)$$

for Taylor- ∞ . Together with an appropriate root point \tilde{x} for which both $f(\tilde{x})$ and ε become negligible, we reach the approximation

$$f(x) \approx \sum_{i=1}^d T_i(x, \tilde{x}). \quad (15)$$

We consequently define the explanation Taylor- ∞ as $E_i(x) = T_i(x, \tilde{x})$. Again, the question of finding suitable root points remains open.

Note that Taylor- ∞ is a strict improvement over Taylor-1 in terms of conservation for degree-1 trigonometric polynomials. A generalization of Taylor- ∞ to higher-degree polynomials is left for future research. Note also that, Taylor- ∞ ultimately being a black-box explanation, it could also be used for other types of functions, but in general the good approximation condition (conservation) cannot be expected to hold in general beyond degree-1 trigonometric polynomials.

B. Quantum layerwise relevance propagation

We start by regrouping the PQC scheme introduced above into a two-step process.

1. Prepare a data-dependent state: $x \mapsto \rho(x)$.
2. Measure a task-dependent observable: $\mathcal{M}(\vartheta) \mapsto \langle \mathcal{M}(\vartheta) \rangle_{\rho(x)}$.

In this *interaction picture*, both the initial state and the measurement observable have been evolved, the former only under data-dependent unitary gates, and the latter only under the trainable ones. This perspective yields the simple computational graph

$$x \xrightarrow{\text{Encoding step}} \rho(x) \xrightarrow{\text{Linear step}} f(x) = \text{tr}\{\rho(x)\mathcal{M}(\vartheta)\}. \quad (16)$$

First the *encoding step* is in general a non-linear map of x , then the *linear step* corresponds to the quantum measurement, and is indeed a linear map of $\rho(x)$. For ease of implementation we introduce a *twin neural network* (twINN) in Appendix D that mirrors these two steps.

Following the divide-and-explain principle underlying the LRP method, we consider a two step explanation, following the computational graph in reverse order. In the first step, we produce the intermediate relevance for $\rho(x)$, which we denote by $R(\rho)$, and which depends on both the value of the function $f(x)$ and the entries of $\rho(x)$ itself: $(f(x), \rho(x)) \mapsto R(\rho)$. In the second step, we reach the final explanation $E(x)$ by using information only from x, ρ , and the relevance of ρ : $(x, \rho, R(\rho)) \mapsto E(x)$. To fully specify this approach one needs only provide concrete formulas for $R(\rho)$ and $E(x)$, which we call the *linear rule* and the *encoding rule*, respectively.

In the XQML method we propose, the linear rule exploits the fact that $f(x)$ is a linear function

$$f(x) = \text{tr}\{\rho(x)\mathcal{M}(\vartheta)\} = \sum_{i,j} \rho_{i,j}(x) \mathcal{M}_{i,j}(\vartheta) \quad (17)$$

of the entries of $\rho(x)$. The linear rule then is simply³ $R_{i,j}(\rho) = \rho_{i,j}(x) \mathcal{M}_{i,j}(\vartheta)$, which is always conservative by construction, and it corresponds to `grad×input` introduced above. The encoding rule is in turn a bit more involved, as not only is the encoding step non-linear, but also this rule must take the relevance $R(\rho)$ into account. For the PQC's we consider (where each input component is encoded only once) each entry of $\rho(x)$ is itself a degree-1 polynomial in each of the components of x . That means that we can consider again the expansion

$$\rho_{i,j}(x) = \rho_{i,j}(\tilde{x}_{i,j}) + \sum_{k=1}^d T_k^{i,j}(x, \tilde{x}_{i,j}) + \varepsilon \quad (18)$$

we have used for the derivation of the Taylor- ∞ relevance. Assuming we find good root points $\tilde{x}_{i,j}$ for each entry of $\rho(x)$, we need to combine these $T_k^{i,j}$ functions with the intermediate

³ For ease of implementation, we actually consider a real-valued version of every complex-valued matrix involved, as we explain in Appendix D.

explanation we obtained from the linear rule. We propose the encoding rule

$$E_k(x) = \sum_{i,j} T_k^{i,j}(x, \tilde{x}_{i,j}) \frac{R_{i,j}(\rho)}{\rho_{i,j}}, \quad (19)$$

with the goal of optimizing for conservation. For this formula only, we use the convention $0/0 = 0$, as for the zero entries $\rho_{i,j}(x) = 0$ it holds also that $R_{i,j}(x) = 0$ from the linear rule, and we need just remove those entries from the sum. One can readily see that, assuming we have good root points, the identity $\sum_{k=1}^d E_k(x) \approx f(x)$ holds, and we call the resulting explanation technique QLRP, which stands for *quantum layerwise relevance propagation*. In Appendix E, we give an algorithm for finding root points.

The QLRP algorithm we propose involves several steps which require exponential space on a classical computer. The number of independent entries of $\rho(x)$ and $\mathcal{M}(\vartheta)$ in general is $\mathcal{O}(\exp(n))$, which accounts for the exponential complexity of the linear rule. The runtime of the algorithms we propose to find the root points for each $\rho_{i,j}(x)$ is polynomial in n , but since there are exponentially many of them, the encoding rule also has exponential complexity. This way our proposed QLRP requires full classical simulation of the quantum circuit and storing the intermediate data-dependent state. In its current formulation, access to the PQC that evaluates the function does not improve the total runtime, as the storage requirements are still exponential throughout.

This explanation technique represents the first adaption of the divide-and-explain principle underlying the LRP explanation technique to QML, albeit with some practical limitations. On the one hand, this is precisely in line with the improvement we wanted to achieve with respect to black-box methods. On the other hand, though, it means that this method is not scalable to larger problems with more qubits involved.

Two important questions remain open. First, it would be interesting to establish how applicable this method is to strictly classically-simulable circuits (like circuits where all states involved are either low-entanglement or low-magic states). Second, one could adapt QLRP to take advantage of quantum hardware which, together with intermediate measurements if necessary, would reduce the storage requirements to be classically tractable, possibly at the cost of only approximately accurate explanations.

IV. NUMERICAL RESULTS

In order to showcase our proposed methods, `Taylor-∞` and QLRP, we devise a proof of concept classification experiment. In this experiment, we use synthetic data for which we can evaluate not only the performance of a PQC, but also for which the explanation quality is easily quantifiable. Our goal is to display the practical application of our protocols, not to claim that quantum models have an advantage over classical algorithms in terms of explainability, nor that our proposed explanation techniques are superior to previously-established ones. With the aim of offering a point of comparison, we

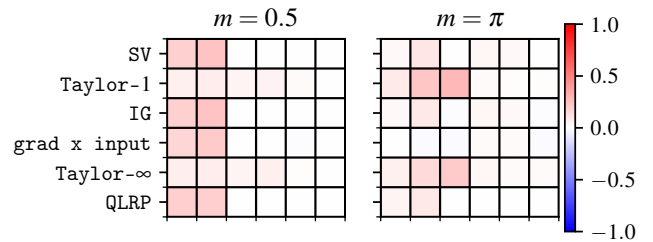


Figure 2. Average predicted explanation for data in the first class, for two synthetic learning tasks. For this class, the relevant components should be the first three. For the right task, with $m = \pi$, we observe that most explanation methods have average relevance close to 0 for all components.

benchmark a subset of the model-agnostic explanations introduced in Section II C together with our novel model-specific techniques.

We propose a classification task using 6-dimensional data sorted into 4 different classes. The specifics of the task can be found in Appendix G. The main feature of the data generation is that three components are relevant for each class. The classification task is rather simple, in that points from different classes cluster together and are close to linearly separable. The relevant components are sampled from a normal distribution centered away from the origin, and the irrelevant components are sampled uniformly random. With this, for each class, we immediately know which components should be given high relevance, and which ones should receive zero attribution.

We train a 6-qubit PQC to solve the multi-class task, where four different fixed observables are measured at the end, one producing a relative score for each class. A relatively simple PQC where each component is encoded only once is already able to solve the task satisfactorily. Then, we produce local explanations for individual input patterns following all the ideas introduced above: both the established ones from Section II C and our novel ones from Section III. In Fig. 2, we can observe the average explanations produced by several different methods for elements of the first class. We note that, although we manually set three components as relevant for each class, from the explainability analysis we can infer that our simple model only makes use of two such components. A similar behavior has been reported in Ref. [21]. We also see a dichotomy between the experiment repetitions: the explanations produced by both `Taylor-1` and `Taylor-∞` differ from the other methods in that the relevant components are more starkly highlighted in the $m = \pi$ repetition.

In order to quantify the quality of explanation, we focus on quality metrics that presuppose known information about the problem. We use a *mask* $M(x)$, which has the same shape as x , and whose entries flag the actually relevant features of the input. Quality evaluation metrics correspond to different similarity measures between the explanation $E(x)$ and the mask $M(x)$ for each given input x . In Appendix F, we provide a precise formulation of the quality metrics considered in this work. In the first measure, called *explana-*

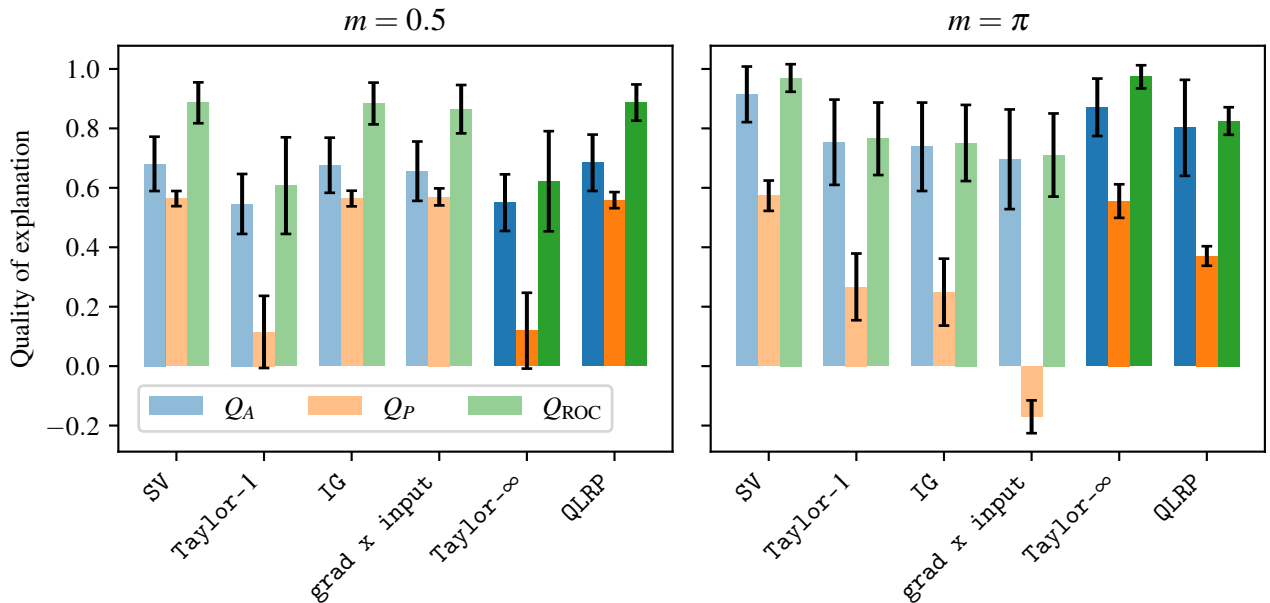


Figure 3. Quality of explanation for methods covered in this work, full report can be found in Appendix G and Fig. 4. The difference in opacity differentiates between the existing local attribution methods introduced in Section II C and the ones we introduce for quantum learning models in Section III.

tion alignment $Q_A(x)$, we simply compute the share of relevance that is distributed to the correct features. This takes the form of a normalized inner product between the element-wise magnitudes of the explanation and the mask. Next, given a set of inputs $(x_i)_i$, another measure is the *Pearson correlation* between the explanations and the masks: $Q_P(x) = \text{Corr}_{\{x_i\}_i}(E(x_i), M(x_i))$. Finally, borrowing from other corners of machine learning, we include a third metric Q_{ROC} based on the *receiver operator characteristics* (ROC), which relies on a binary classification of explanations based on the mask. The quality of explanation is depicted in Fig. 3.

As explained in Appendix G, we perform three different experiment repetitions, labeled by a hyperparameter we call m . We consider three values of m , and we observe that the larger m is, the more difficult the classification becomes. In Fig. 3 we show only some of the explanation methods for the intermediate and higher values of the hyperparameter: $m = 0.5$ and $m = \pi$ respectively. In both cases the achieved test accuracy is above 80%, and in the intermediate case it is close to optimal. What we observe in the plots is a rich range of qualities, both for different methods and for different metrics. Without a clear best choice across the board, the messages to be extracted from these experiments are:

- There is no obvious benefit in using model-specific explanations over model-agnostic ones. This may be an intrinsic fact of QXML or an artifact of the simplicity of our synthetic classification task.
- If the goal is to quantify the explainability of a QML model for a specific data-set, then efforts must be invested in validating each QXML method for that set.

- If the goal is to quantify the explainability of a QML model in general, then further research is needed. In this sense, there is no free lunch also for explainability.

These messages are in line with what we see in classical ML, where explanation methods start distinguishing themselves only with increasing complexity of the model and difficulty of the task. For instance, for very deep convolutional NNs or for transformers gradient-based approaches start having problems (while they may work fine for simple settings) [88].

When comparing the left and right plots in Fig. 3, we observe that most explanations achieve higher scores for both the explanation alignment Q_A and the ROC-based Q_{ROC} evaluation metrics for the $m = \pi$ task. For the XQML methods introduced in this work, we see a reverse in the performance ranking: QLRP outperforms Taylor- ∞ for $m = 0.5$, but the converse is true for $m = \pi$. This may well be from specific artifacts in the particular toy problem we employ. Structurally, both experiment repetitions are very similar, with the only noticeable difference that the simple classifier achieves better accuracy for $m = 0.5$ than for $m = \pi$. From this, we hypothesize that the second task being more complex, the training algorithm settles for a slightly simpler model than in the first task, thus accounting for the lower accuracy. If this were the case, a possible explanation for the better performance of the gradient-based explanation methods could be that the second model is closer to a linear one. As put forth in Ref. [88] for the study of explainability in transformers, linearity could account for a generic improvement in the explanation quality of gradient-based methods. Nevertheless, this particular experiment is synthetic in nature, and there is no strong evidence

to believe the trends showed here would generalize to larger, more realistic scenarios.

Next to measuring the quality of explanation across several metrics, we recall that the main difference between Taylor-1 and $\text{Taylor-}\infty$ ought to be an improvement in conservation. To assess potential differences, we consider the *relative approximation error*

$$\frac{|\sum_{i=1}^d E_i(x) - f(x)|}{|f(x)|} \quad (20)$$

of a specific explanation E . In Table III we report the relative approximation error of both Taylor-1 and $\text{Taylor-}\infty$ averaged over the data-sets corresponding to each of the two experiments.

Table III. Relative approximation errors achieved by Taylor-1 and $\text{Taylor-}\infty$ on the synthetic classification tasks.

XQML Method	$m = 0.5$	$m = \pi$
Taylor-1	0.032 ± 0.005	0.40 ± 0.03
$\text{Taylor-}\infty$	0.028 ± 0.004	0.52 ± 0.04

We briefly recall that one of the guiding principles for the derivation of $\text{Taylor-}\infty$ was to include more terms in the Taylor expansion of the function. By exploiting the structure of low-degree trigonometric functions, one would have expected a high degree of conservation in the resulting explanations. This is not what we observe in these experiments, where both methods produce similar results, and Taylor-1 is even more conservative for the $m = \pi$ experiment. The hypothesis introduced above that the classifier in $m = \pi$ was closer to linear would also explain this disparity. While $\text{Taylor-}\infty$ contains strictly more terms of the Taylor expansion than Taylor-1 , it could be that the inclusion of these extra terms is harmful if the underlying function is close to linear, as for example if only positive contributions are taken from a balanced sum. Again, we leave this as a conjecture, and we do not spend further efforts in characterizing behaviors that could be specific to this synthetic toy problem.

V. SUMMARY AND OUTLOOK

In this work, we have charted what the field of *explainable quantum machine learning* (XQML) could look like, as a “preview” of the field. In order to provide such a preview, we have discussed the importance of explainability in general for machine learning systems of all kinds. After reviewing the main concepts from the fields of *explainable artificial intelligence* (XAI) and *quantum machine learning* (QML), we have noted there is one significant disparity between QML and classical ML: while XAI is substantially more developed than XQML, advancements in explainability have lagged behind breakthroughs in performance for deep learning models. This is not yet the case for QML models, and this offers an opportunity for avoiding a trade-off between explainability and performance in QML. By recognizing the depth of potential

impact of explainability in the development of QML, we can, therefore, decide to prioritize the design of inherently explainable QML models.

Along the way, we have remarked that the very core features of quantum mechanics may render impossible the direct application of certain XAI ideas to QML models. The hardness of efficiently representing quantum states classically, the impossibility of copying quantum information, and the destruction of quantum coherence when measuring intermediate steps of a quantum computation force us to be imaginative in the design of XQML techniques. Next to laying out the formalism of XQML, in this work we have also taken the first steps forward, by proposing quantum versions of well-established classical XAI techniques.

In our work, we have focused on post-hoc local attribution methods on the one hand, and on QML models based on *parametrized quantum circuits* (PQCs) on the other hand. We have introduced two novel explanation techniques to be applicable for PQC-based QML models. Referencing the classical XAI ideas from where they stem, we called them $\text{Taylor-}\infty$ and *quantum layerwise relevance propagation*, QLRP. $\text{Taylor-}\infty$ leverages the Taylor decomposition of the labeling function, and allows us to retain a higher-degree of information than its classical counterpart Taylor-1 . $\text{Taylor-}\infty$ is a black-box explanation, but it exploits the trigonometric structure of usual QML labeling functions [27]. QLRP adapts the divide-and-explain mindset adopted in *layerwise relevance propagation* (LRP) [52] to PQC-based QML models. By considering a PQC as a linear model on a high-dimensional Hilbert space, in QLRP we propose a two-step explanation, taking into account the structure of the intermediate quantum feature map.

We present these novel explanation techniques with derivations substantially expanding on their classical counterparts. We accompany these analytical results with numerical experiments using a synthetic learning task. The goal of these experiments is not so much to present state-of-the-art simulations at the forefront of classical simulations, but have been set up mainly to transparently illustrate the pipeline of implementation and evaluation of XQML techniques.

By reporting the performance of different explanation techniques, we provide a comparison between existing model-agnostic explanations and our novel model-specific ones. We leave for future work a full characterization of the potential advantages and merits of the different techniques. We decide on this humbler stance for two important reasons: First, it is our understanding that strong evidence for QML models that perform well on practically-relevant tasks is still missing [85], which means there is a lack of understanding of which features will prove eventually relevant in the design of QML models. Second, the large-scale benchmark from Ref. [89] raised the alarm that cherry-picking is a common issue in QML research. The conclusion from our numerics is that different methods perform differently for different tasks and under different metrics, and so further research is needed in (1) characterizing all XQML techniques further, and (2) understanding which features and metrics will be relevant when eventually designing practically-relevant QML models. To be

precise: we do not search for quantum advantage in XAI, we rather contribute to the foundations of XQML and take first meaningful steps.

Arguably the main challenge we encountered is that of scalability. To design an explanation that makes use of the specific structure of a PQC, the straightforward approach of fully characterizing the Hilbert space of quantum computation quickly becomes classically intractable. A generic issue in QML also arises in XQML: the choice of encoding is critical [27]. Namely, independently of the achieved performance in solving a learning task, if the chosen encoding introduces too much complexity, explanation becomes critically difficult. We identified a trade-off that is also present in XAI for deep learning: the less structured the learning model, the more complex the explanation method must be.

The challenges we encountered additionally dictate the potential ways forward in XQML. The scalability issue must be addressed by designing smart quantum explanation algorithms that are quantum-time- and classical-space-efficient. Such algorithms would involve intermediate measurements to extract partial information about the quantum computation. To better grasp the effect of the encoding on explainability, a promising direction is studying the explainability of kernel-based learning models, since QML has a deep connection to kernel methods [90]. Additionally, the structure-complexity trade-off could dictate the design of QML models: to have the specific structures that would make explainability easier (for example, as proposed in Ref. [91]).

Other promising research directions include the generalization of our explanation techniques to other encodings, including data re-uploading circuits [28], or QML models taking quantum states as input data. Implementing larger-scale experiments would help us towards understanding what desirable features arise from different explanation techniques and evaluation metrics. From both these considerations, a further impactful contribution would be the compilation of earnest

explainability-based guidelines to the design of PQCs. With these guidelines, the difficulty of finding specific application domains for QML could be partially alleviated.

More generally, other bridges could be built between XAI and quantum computing. For instance, one could further study the possibility of explaining parts of the quantum circuit, instead of only the produced explanations for a learning task [23]. This direction includes questions around mechanistic interpretability, meaning the study of relevance at the level of individual neurons. In the opposite direction, an interesting contribution would be to better chart the space of potential quantum advantage in XAI by establishing the computational complexity of the different steps in the current XAI pipelines for deep learning. Finally, our work could be expanded toward the direction of hybrid quantum-classical ML models, where PQCs are directly combined with classical neural networks.

Acknowledgements

The authors thank Anna Dawid for insightful comments in an earlier version of this draft. This work has been supported by the BMBF (Hybrid), the Munich Quantum Valley (K-4 and K-8), the BMWK (EniQmA), the Quantum Flagship (PasQuans2, Millenion), QuantERA (HQCC), the Cluster of Excellence MATH+, the DFG (CRC 183 and Research Unit KI-FOR 5363 – project ID: 459422098), the Einstein Foundation (Einstein Research Unit on Quantum Devices), Berlin Quantum, and the ERC (DebuQC). For the Einstein Research Unit, this has been the result of an important joint-node collaboration. EGF is supported by a Google PhD Fellowship and acknowledges travel support from the European Union’s Horizon 2020 research and innovation programme under grant agreement No 951847.

-
- [1] Y. A. LeCun, L. Bottou, G. B. Orr, and K.-R. Müller, in *Neural networks: Tricks of the trade* (Springer, Berlin, 2012) pp. 9–48.
- [2] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning* (MIT Press, Cambridge (MA), 2016).
- [3] O. Nov, N. Singh, and D. Mann, Putting ChatGPT’s medical advice to the (Turing) test: survey study, *JMIR Med. Educ.* **9**, e46939 (2023).
- [4] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. V. D. Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot, S. Dieleman, D. Grewe, J. Nham, N. Kalchbrenner, I. Sutskever, T. Graepel, T. Lillicrap, M. Leach, K. Kavukcuoglu, and D. Hassabis, Mastering the game of go with deep neural networks and tree search, *Nature* **529**, 484 (2016).
- [5] L. Deng, G. Hinton, and B. Kingsbury, New types of deep neural network learning for speech recognition and related applications: An overview, *IEEE Int. Conf. Ac. Speech and Signal Proc. (ICASSP)*, 8599 (2013).
- [6] K. He, X. Zhang, S. Ren, and J. Sun, Deep residual learning for image recognition, *Proc. IEEE Conf. Comp. Vision and Pattern Rec. (CVPR)*, 770 (2016).
- [7] F. Arute *et al.*, Quantum supremacy using a programmable superconducting processor, *Nature* **574**, 505 (2019).
- [8] D. Hangleiter and J. Eisert, Computational advantage of quantum random sampling, *Rev. Mod. Phys.* **95**, 035001 (2023).
- [9] J. Biamonte, P. Wittek, N. Pancotti, P. Rebentrost, N. Wiebe, and S. Lloyd, Quantum machine learning, *Nature* **549**, 195 (2017).
- [10] V. Dunjko and H. J. Briegel, Machine learning & artificial intelligence in the quantum domain: a review of recent progress, *Rep. Prog. Phys.* **81**, 074001 (2018).
- [11] G. Carleo, I. Cirac, K. Cranmer, L. Daudet, M. Schuld, N. Tishby, L. Vogt-Maranto, and L. Zdeborová, Machine learning and the physical sciences, *Rev. Mod. Phys.* **91**, 045002 (2019).
- [12] J. R. McClean, J. Romero, R. Babbush, and A. Aspuru-Guzik, The theory of variational hybrid quantum-classical algorithms, *New J. Phys.* **18**, 023023 (2016).
- [13] R. Sweke, J.-P. Seifert, D. Hangleiter, and J. Eisert, On the quantum versus classical learnability of discrete distributions,

- Quantum* **5**, 417 (2021).
- [14] Y. Liu, S. Arunachalam, and K. Temme, A rigorous and robust quantum speed-up in supervised machine learning, *Nature Phys.* **17**, 1013 (2021).
- [15] N. Pirnay, R. Sweke, J. Eisert, and J.-P. Seifert, A superpolynomial quantum-classical separation for density modelling, *Phys. Rev. A* **107**, 042416 (2023).
- [16] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. J. Goodfellow, and R. Fergus, Intriguing properties of neural networks, *ICLR (Poster)* (2014).
- [17] S. Lapuschkin, A. Binder, G. Montavon, K.-R. Müller, and W. Samek, Analyzing classifiers: Fisher vectors and deep neural networks, in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016) pp. 2912–2920.
- [18] W. Samek, G. Montavon, S. Lapuschkin, C. J. Anders, and K.-R. Müller, Explaining deep neural networks and beyond: A review of methods and applications, *Proc. IEEE* **109**, 247 (2021).
- [19] D. Minh, H. Wang, Y. Li, and T. N. Nguyen, Explainable artificial intelligence: a comprehensive review, *Artif. Intell. Rev.* **55**, 3503 (2022).
- [20] A. B. Arrieta, N. D. Rodríguez, J. D. Ser, A. Bennetot, S. Tabik, A. Barbado, S. García, S. Gil-Lopez, D. Molina, R. Benjamins, R. Chatila, and F. Herrera, Explainable artificial intelligence (XAI): concepts, taxonomies, opportunities and challenges toward responsible AI, *Inf. Fusion* **58**, 82 (2020).
- [21] P. Steinmüller, T. Schulz, F. Graf, and D. Herr, eXplainable AI for quantum machine learning, [arXiv:2211.01441](https://arxiv.org/abs/2211.01441) (2022).
- [22] R. Lifshitz, Quantum deep dreaming: A novel approach for quantum circuit design, [arXiv:2211.04343](https://arxiv.org/abs/2211.04343) (2022).
- [23] R. Heese, T. Gerlach, S. Mücke, S. Müller, M. Jakobs, and N. Piatkowski, Explaining quantum circuits with shapley values: Towards explainable quantum machine learning, [arXiv:2301.09138](https://arxiv.org/abs/2301.09138) (2023).
- [24] L. Pira and C. Ferrie, On the interpretability of quantum neural networks, *Quant. Mach. Int.* **6**, 52 (2024).
- [25] G. K. Dziugaite, S. Ben-David, and D. M. Roy, Enforcing interpretability and its statistical impacts: Trade-offs between accuracy and interpretability, [arXiv:2010.13764](https://arxiv.org/abs/2010.13764) (2020).
- [26] F. J. Gil Vidal and D. O. Theis, Input redundancy for parameterized quantum circuits, *Frontiers in Physics* **8**, 13 (2020).
- [27] M. Schuld, R. Sweke, and J. J. Meyer, Effect of data encoding on the expressive power of variational quantum-machine-learning models, *Phys. Rev. A* **103**, 032430 (2021).
- [28] A. Pérez-Salinas, A. Cervera-Lierta, E. Gil-Fuster, and J. I. Latorre, Data re-uploading for a universal quantum classifier, *Quantum* **4**, 226 (2020).
- [29] A. Abbas, R. King, H.-Y. Huang, W. J. Huggins, R. Movassagh, D. Gilboa, and J. McClean, On quantum backpropagation, information reuse, and cheating measurement collapse, *Adv. Neur. Inf. Proc. Syst.* **2026**, 36 (2024).
- [30] J. R. McClean, S. Boixo, V. N. Smelyanskiy, R. Babbush, and H. Neven, Barren plateaus in quantum neural network training landscapes, *Nature Comm.* **9**, 4812 (2018).
- [31] M. Larocca, S. Thanasilp, S. Wang, K. Sharma, J. Biamonte, P. J. Coles, L. Cincio, J. R. McClean, Z. Holmes, and M. Cerezo, A review of barren plateaus in variational quantum computing, [arXiv:2405.00781](https://arxiv.org/abs/2405.00781) (2024).
- [32] S. Bravyi and D. Gosset, Improved classical simulation of quantum circuits dominated by Clifford gates, *Phys. Rev. Lett.* **116**, 250501 (2016).
- [33] E. Tang, A quantum-inspired classical algorithm for recommendation systems, in *Proceedings of the 51st Annual ACM SIGACT Symposium on Theory of Computing*, STOC '19 (ACM, 2019).
- [34] F. J. Schreiber, J. Eisert, and J. J. Meyer, Classical surrogates for quantum learning models, *Phys. Rev. Lett.* **131**, 100803 (2023).
- [35] J. Landman, S. Thabet, C. Dalyac, H. Mhiri, and E. Kashefi, Classically approximating variational quantum machine learning with random fourier features, in *The Eleventh International Conference on Learning Representations* (2023).
- [36] R. Sweke, E. Recio, S. Jerbi, E. Gil-Fuster, B. Fuller, J. Eisert, and J. J. Meyer, Potential and limitations of random Fourier features for dequantizing quantum machine learning, [arXiv:2309.11647](https://arxiv.org/abs/2309.11647) (2023).
- [37] S. Shin, Y. S. Teo, and H. Jeong, Dequantizing quantum machine learning models using tensor networks, *Phys. Rev. Res.* **6**, 023218 (2024).
- [38] M. S. Rudolph, E. Fontana, Z. Holmes, and L. Cincio, Classical surrogate simulation of quantum systems with LOWESA, [arXiv:2308.09109](https://arxiv.org/abs/2308.09109) (2023).
- [39] M. Cerezo, M. Larocca, D. García-Martín, N. L. Diaz, P. Braccia, E. Fontana, M. S. Rudolph, P. Bermejo, A. Ijaz, S. Thanasilp, *et al.*, Does provable absence of barren plateaus imply classical simulability? or, why we need to rethink variational quantum computing, [arXiv:2312.09121](https://arxiv.org/abs/2312.09121) (2023).
- [40] B. Dias and R. Koenig, Classical simulation of non-Gaussian fermionic circuits, *Quantum* **8**, 1350 (2024).
- [41] D. Gunning and D. W. Aha, DARPA's explainable artificial intelligence (XAI) program, *AI Mag.* **40**, 44 (2019).
- [42] S. Lapuschkin, S. Wäldchen, A. Binder, G. Montavon, W. Samek, and K.-R. Müller, Unmasking clever Hans predictors and assessing what machines really learn, *Nature Comm.* **10**, 1096 (2019).
- [43] J. R. Kauffmann, J. Dippel, L. Ruff, W. Samek, K. Müller, and G. Montavon, The clever Hans effect in unsupervised learning, [arXiv:2408.08041](https://arxiv.org/abs/2408.08041) (2024).
- [44] R. Geirhos, J.-H. Jacobsen, C. Michaelis, R. Zemel, W. Brendel, M. Bethge, and F. A. Wichmann, Shortcut learning in deep neural networks, *Nature Mach. Intell.* **2**, 665–673 (2020).
- [45] A. J. DeGrave, J. D. Janizek, and S. Lee, AI for radiographic COVID-19 detection selects shortcuts over signal, *Nature Mach. Intell.* **3**, 610 (2021).
- [46] A. S. Ross, M. C. Hughes, and F. Doshi-Velez, Right for the right reasons: Training differentiable models by constraining their explanations, in *IJCAI* (ijcai.org, 2017) pp. 2662–2670.
- [47] C. J. Anders, L. Weber, D. Neumann, W. Samek, K. Müller, and S. Lapuschkin, Finding and removing clever Hans: Using explanation methods to debug and improve deep models, *Inf. Fusion* **77**, 261 (2022).
- [48] P. Chormai, J. Herrmann, K. Müller, and G. Montavon, Disentangled explanations of neural network predictions by finding relevant subspaces, *IEEE Trans. Pattern Anal. Mach. Intell.* **46**, 7283 (2024).
- [49] M. Krenn, R. Pollice, S. Y. Guo, M. Aldeghi, A. Cervera-Lierta, P. Friederich, G. dos Passos Gomes, F. Häse, A. Jinich, A. Nigam, Z. Yao, and A. Aspuru-Guzik, On scientific understanding with artificial intelligence, *Nature Rev. Phys.* **4**, 761–769 (2022).
- [50] P. Schramowski, W. Stammer, S. Teso, A. Brugger, F. Herbert, X. Shao, H. Luigs, A. Mahlein, and K. Kersting, Making deep neural networks right for the right scientific reasons by interacting with their explanations, *Nature Mach. Intell.* **2**, 476 (2020).
- [51] M. Dreyer, R. Achibat, W. Samek, and S. Lapuschkin, Understanding the (extra-)ordinary: Validating deep model decisions with prototypical concept-based explanations, in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops* (2024) pp. 3491–3501.

- [52] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, and W. Samek, On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation, *PloS one* **10**, e0130140 (2015).
- [53] T. Schnake, O. Eberle, J. Lederer, S. Nakajima, K. T. Schutt, K.-R. Müller, and G. Montavon, Higher-order explanations of graph neural networks via relevant walks, *IEEE Trans. Pattern Anal. Mach. Intell.* **44**, 7581 (2022).
- [54] B. Kim, M. Wattenberg, J. Gilmer, C. J. Cai, J. Wexler, F. B. Viégas, and R. Sayres, Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV), in *ICML*, Proceedings of Machine Learning Research, Vol. 80 (PMLR, 2018) pp. 2673–2682.
- [55] R. Achibat, M. Dreyer, I. Eisenbraun, S. Bosse, T. Wiegand, W. Samek, and S. Lapuschkin, From attribution maps to human-understandable explanations through concept relevance propagation, *Nature Mach. Intell.* **5**, 1006 (2023).
- [56] J. Vielhaben, S. Lapuschkin, G. Montavon, and W. Samek, Explainable AI for time series via virtual inspection layers, *Pattern Recognit.* **150**, 110309 (2024).
- [57] Y. Lou, R. Caruana, and J. Gehrke, Intelligible models for classification and regression, in *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '12 (Association for Computing Machinery, New York, NY, USA, 2012) p. 150–158.
- [58] B. Zhou, A. Khosla, À. Lapedriza, A. Oliva, and A. Torralba, Learning deep features for discriminative localization, in *CVPR* (IEEE Computer Society, 2016) pp. 2921–2929.
- [59] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, Attention is all you need, in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, NIPS'17 (Curran Associates Inc., Red Hook, NY, USA, 2017) p. 6000–6010.
- [60] T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. Duvenaud, Neural ordinary differential equations, in *NeurIPS* (2018) pp. 6572–6583.
- [61] M. Sundararajan, A. Taly, and Q. Yan, Axiomatic attribution for deep networks, in *International conference on machine learning* (PMLR, 2017) pp. 3319–3328.
- [62] E. Štrumbelj and I. Kononenko, An efficient explanation of individual classifications using game theory, *J. Mach. Learn. Res.* **11**, 1 (2010).
- [63] A.-K. Dombrowski, J. E. Gerken, K.-R. Müller, and P. Kessel, Diffeomorphic counterfactuals with generative models, *IEEE Trans. Patt. Ana. Mach. Int.* (2023).
- [64] A. Nguyen, A. Dosovitskiy, J. Yosinski, T. Brox, and J. Clune, Synthesizing the preferred inputs for neurons in neural networks via deep generator networks, *Adv. Neu. Inf. Proc. Sys.* (NeurIPS) **29**, 3395 (2016).
- [65] Google Research, Inceptionism: Going deeper into neural networks, <https://ai.googleblog.com/2015/06/inceptionism-going-deeper-into-neural.html> (2015), accessed: 2024-08-01.
- [66] L. S. Shapley, A value for n -person games, in *Contributions to the Theory of Games (AM-28), Volume II* (Princeton University Press, 1953) pp. 307–317.
- [67] S. M. Lundberg and S. Lee, A unified approach to interpreting model predictions, *NIPS*, 4765 (2017).
- [68] M. Sundararajan and A. Najmi, The many shapley values for model explanation, in *International conference on machine learning* (PMLR, 2020) pp. 9269–9278.
- [69] J. Castro, D. Gómez, and J. Tejada, Polynomial calculation of the Shapley value based on sampling, *Computers & Operations Research Selected Papers Presented at the Tenth International Symposium on Locational Decisions (ISOLDE X)*, **36**, 1726 (2009).
- [70] M. Ancona, C. Oztireli, and M. Gross, Explaining deep neural networks with a polynomial time algorithm for shapley value approximation, in *International Conference on Machine Learning* (PMLR, 2019) pp. 272–281.
- [71] D. E. Rumelhart, G. E. Hinton, and R. J. Williams, Learning representations by back-propagating errors, *Nature* **323**, 533–536 (1986).
- [72] G. Montavon, J. R. Kauffmann, W. Samek, and K. Müller, Explaining the predictions of unsupervised learning models, in *xxAI@ICML*, Lecture Notes in Computer Science, Vol. 13200 (Springer, 2020) pp. 117–138.
- [73] A. Ali, T. Schnake, O. Eberle, G. Montavon, K. Müller, and L. Wolf, XAI for transformers: Better explanations through conservative propagation, in *ICML*, Proceedings of Machine Learning Research, Vol. 162 (PMLR, 2022) pp. 435–451.
- [74] G. Montavon, A. Binder, S. Lapuschkin, W. Samek, and K.-R. Müller, Layer-wise relevance propagation: An overview, in *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*, Vol. 11700, edited by W. Samek, G. Montavon, A. Vedaldi, L. K. Hansen, and K.-R. Müller (Springer International Publishing, 2019) pp. 193–209.
- [75] W. Samek, A. Binder, G. Montavon, S. Lapuschkin, and K. Müller, Evaluating the visualization of what a deep neural network has learned, *IEEE Trans. Neural Networks Learn. Syst.* **28**, 2660 (2017).
- [76] L. Arras, A. Osman, and W. Samek, CLEVR-XAI: A benchmark dataset for the ground truth evaluation of neural network explanations, *Information Fusion* **81**, 14–40 (2022).
- [77] A. Hedström, L. Weber, D. Bareeva, D. Krakowczyk, F. Motzkus, W. Samek, S. Lapuschkin, and M. M.-C. Höhne, Quantus: An explainable ai toolkit for responsible evaluation of neural network explanations and beyond, *J. Mach. Learn. Res.* **24**, 1 (2023).
- [78] W. R. Swartout and J. D. Moore, Explanation in second generation expert systems, in *Second generation expert systems* (Springer, 1993) pp. 543–585.
- [79] H.-Y. Huang, R. Kueng, and J. Preskill, Predicting many properties of a quantum system from very few measurements, *Nature Phys.* **16**, 1050 (2020).
- [80] J. I. Cirac, D. Pérez-García, N. Schuch, and F. Verstraete, Matrix product states and projected entangled pair states: Concepts, symmetries, theorems, *Rev. Mod. Phys.* **93**, 045003 (2021).
- [81] J. Eisert, M. Cramer, and M. B. Plenio, Area laws for the entanglement entropy, *Rev. Mod. Phys.* **82**, 277 (2010).
- [82] D. Gottesman, Theory of fault-tolerant quantum computation, *Phys. Rev. A* **57**, 127 (1998).
- [83] S. Masot-Llima and A. Garcia-Saez, Stabilizer tensor networks: universal quantum simulator on a basis of stabilizer states, *arXiv:2403.08724* (2024).
- [84] K. Chinzei, S. Yamano, Q. H. Tran, Y. Endo, and H. Oshima, Trade-off between gradient measurement efficiency and expressivity in deep quantum neural networks, *arXiv:2406.18316* (2024).
- [85] E. Gil-Fuster, C. Gyurik, A. Pérez-Salinas, and V. Dunjko, On the relation between trainability and dequantization of variational quantum learning models, *arXiv:2406.07072* (2024).
- [86] M. T. Ribeiro, S. Singh, and C. Guestrin, "Why should I trust you?": Explaining the predictions of any classifier, in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '16 (Association for Computing Machinery, New York, NY, USA, 2016) p.

- 1135–1144.
- [87] G. Montavon, W. Samek, and K.-R. Müller, Methods for interpreting and understanding deep neural networks, *Dig. Sign. Proc.* **73**, 1 (2018).
- [88] R. Achibat, S. M. V. Hatefi, M. Dreyer, A. Jain, T. Wiegand, S. Lapuschkin, and W. Samek, AttnLRP: Attention-aware layer-wise relevance propagation for transformers, in *Proceedings of the 41st International Conference on Machine Learning*, Proceedings of Machine Learning Research, Vol. 235 (2024) pp. 135–168.
- [89] J. Bowles, S. Ahmed, and M. Schuld, Better than classical? the subtle art of benchmarking quantum machine learning models, [arXiv:2403.07059](https://arxiv.org/abs/2403.07059) (2024).
- [90] M. Schuld, Supervised quantum machine learning models are kernel methods, [arXiv:2101.11020](https://arxiv.org/abs/2101.11020) (2021).
- [91] K. Cybiński, J. Enouen, A. Georges, and A. Dawid, Speak so a physicist can understand you! tetrisenn for detecting phase transitions and order parameters, [arXiv: 2411.02237](https://arxiv.org/abs/2411.02237) (2024).
- [92] G. Montavon, S. Lapuschkin, A. Binder, W. Samek, and K.-R. Müller, Explaining non-linear classification decisions with deep Taylor decomposition, *Pattern Recognition* **65**, 211 (2017).
- [93] M. Schuld, A. Bocharov, K. M. Svore, and N. Wiebe, Circuit-centric quantum classifiers, *Phys. Rev. A* **101**, 032308 (2020).
- [94] V. Bergholm, J. Izaac, M. Schuld, C. Gogolin, S. Ahmed, V. Ajith, M. S. Alam, G. Alonso-Linaje, B. AkashNarayanan, A. Asadi, J. M. Arrazola, U. Azad, S. Banning, C. Blank, T. R. Bromley, B. A. Cordier, J. Ceroni, A. Delgado, O. D. Matteo, A. Dusko, T. Garg, D. Guala, A. Hayes, R. Hill, A. Ijaz, T. Isaacsson, D. Ittah, S. Jahangiri, P. Jain, E. Jiang, A. Khandelwal, K. Kottmann, R. A. Lang, C. Lee, T. Loke, A. Lowe, K. McKiernan, J. J. Meyer, J. A. Montañez-Barrera, R. Moyard, Z. Niu, L. J. O’Riordan, S. Oud, A. Panigrahi, C.-Y. Park, D. Polatajko, N. Quesada, C. Roberts, N. Sá, I. Schoch, B. Shi, S. Shu, S. Sim, A. Singh, I. Strandberg, J. Soni, A. Száva, S. Thabet, R. A. Vargas-Hernández, T. Vincent, N. Vitucci, M. Weber, D. Wierichs, R. Wiersema, M. Willmann, V. Wong, S. Zhang, and N. Killoran, Pennylane: Automatic differentiation of hybrid quantum-classical computations, [arXiv:1811.04968](https://arxiv.org/abs/1811.04968) (2022).
- [95] R. Frostig, M. Johnson, and C. Leary, *Compiling machine learning programs via high-level tracing* (2018).
- [96] D. Smilkov, N. Thorat, B. Kim, F. Viégas, and M. Wattenberg, Smoothgrad: removing noise by adding noise, [arXiv:1706.03825](https://arxiv.org/abs/1706.03825) (2017).

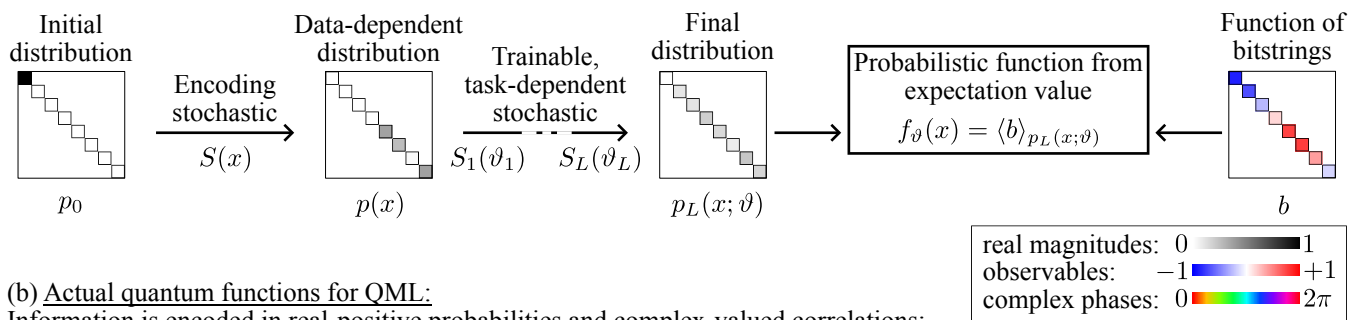
Supplementary Material for “Opportunities and limitations of explaining quantum machine learning”

Appendix A: Probabilistic functions for machine learning

For didactic purposes, one could consider the closest fully-classical analogous to the hypothesis families of quantum functions introduced in Section II A, as shown in Fig. 1. Such an exercise would involve, e.g., specifying a parametrization of the set of doubly-stochastic matrices. A real-valued function could then arise from a probabilistic classical circuit as a three-step process: (1) start from an easy distribution (like the Kronecker delta on the all-0 string), (2) apply a sequence of parametrized doubly-stochastic gates (some of which would be data dependent, and some not), and (3) estimate the expectation value of a fixed function on bits b . Doubly-stochastic transformations play the role of unitary matrices in this hypothetical scenario. Unitary matrices can safely be thought of as the generalization of rotations to complex-valued fields. This way, the action of a unitary on the matrix D_p is akin to rotating the vector p in a way that does not incur issues with the signs of the entries or the normalization. The main difference between doubly stochastic transformations and unitary transformations is that the former can only increase the Shannon entropy of the distribution (the inverse of a doubly stochastic matrix need not be stochastic), whereas the latter also allow for decreasing Shannon entropy, since the inverse of a unitary matrix is also unitary.

One immediately notices that even though one could define classical ML models based on these probabilistic circuits, these are really far from what practitioners use. With this presentation we hope to illustrate the conceptual simplicity of PQC-based QML models, yet in practice the function families PQC give rise to are rich and intricate.

(a) Hypothetical probabilistic functions for ML: Information is encoded in real positive probabilities:



(b) Actual quantum functions for QML:

Information is encoded in real-positive probabilities and complex-valued correlations:

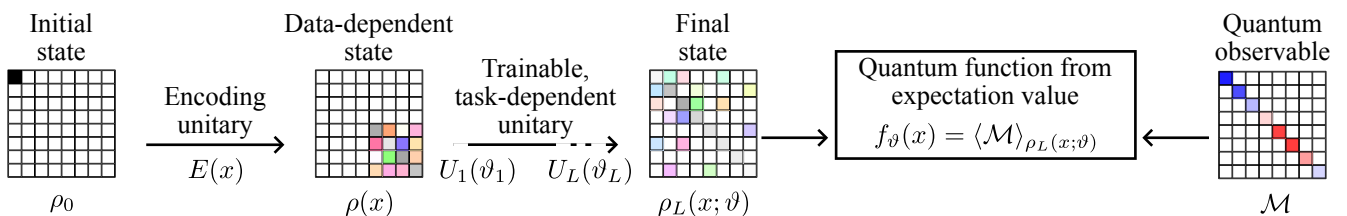


Figure 1. Sketch illustrating the relation between the quantum functions used in QML and the conceptually closest fully-classical probabilistic functions one could use for ML. (a) Classical probability distributions can be represented as diagonal matrices. Then, information on the data or trainable parameters can be encoded as a bi-stochastic transformation, resulting in a parametrized final distribution. A real-valued function can be recovered from evaluating the expectation value of any function of bitstrings with respect to the classical final distribution. (b) The same idea applies, just that now the classical information is stored in quantum states, represented as positive semi-definite, unit trace, Hermitian matrices. The core mechanisms are left unchanged, just now we have access to complex-valued entries outside the main diagonal. Again, a real-valued function can be recovered as the expectation value of a fixed quantum observable with respect to the final parametrized quantum state.

Appendix B: Illustrative presentation of LRP

As introduced in Section II C 0c, we consider functions following a simple computational graph

$$x = z_0 \xrightarrow{\phi_1} z_1 \xrightarrow{\phi_2} \dots \xrightarrow{\phi_l} z_l \xrightarrow{\phi_{l+1}} \dots \xrightarrow{\phi_L} z_L = f(x).$$

Then, the first step in LRP is to produce a relevance heat-map for the penultimate intermediate state, which we denote $R^{(L-1)}(x)$. This intermediate explanation has the same size and shape as z_{L-1} , so it ultimately depends on the structure of the model. In this line, we would say $R^{(L)}(x) = f(x)$, which for simplicity we take to be a single real number. We would write the explanation $R^{(L-1)}(x)$ as a function of $f(x) = z_L$ and z_{L-1} :

$$\begin{array}{ccc} f(x) = z_L = R^{(L)}(x) & \longmapsto & R^{(L-1)}(x). \\ \uparrow \phi_L & \nearrow & \\ z_{L-1} & & \end{array}$$

For ease of notation, from now on we drop the x dependence of the intermediate explanations, but in what follows we talk about explaining the model prediction for a single input (we are studying local explanations).

Now we have an explanation for the penultimate hidden layer, and our goal is to, from here, produce an explanation for the second-to-last one. For the penultimate layer, though, in general we have $z_{L-1} \neq R^{(L-1)}$, which was the case for the last layer. This means that, in order to have an intermediate explanation for the next layer $R^{(L-2)}$, we must also take $R^{(L-1)}$ into account:

$$\begin{array}{ccccc} & & R^{(L-1)} & \longmapsto & R^{(L-2)} \\ & \nearrow & \uparrow & \nearrow & \uparrow \\ R^{(L)} & & z_{L-1} & \xleftarrow{\phi_{L-1}} & z_{L-1} \\ & \searrow \phi_L & & & \end{array}$$

This idea gives us a blueprint that can be followed throughout: the l^{th} intermediate explanation is a simple function of z_l, z_{l+1} and the $(l+1)^{\text{th}}$ intermediate explanation [92]:

$$\begin{array}{ccccccc} & & R^{(L-1)} & \longmapsto & \dots & \longmapsto & R^{(1)} & \longmapsto & R^{(0)} \\ & \nearrow & \uparrow & \nearrow & & \nearrow & \uparrow & \nearrow & \uparrow \\ z_L = R^{(L)} & & z_{L-1} & \xleftarrow{\phi_{L-1}} & \dots & \xleftarrow{\phi_2} & z_1 & \xleftarrow{\phi_1} & z_0 \\ & \searrow \phi_L & & & & & & & \end{array}$$

At this point we recall that $z_0 = x$ and it follows that $R^{(0)}$ is the last explanation, thus the one we produce for this model $E(x) = R^{(0)}$.

Appendix C: Derivation of the Taylor- ∞ explanation

Proposition 1 (Taylor explanation). *Let $f: \mathbb{R}^d \rightarrow \mathbb{R}$ be a trigonometric polynomial*

$$f(x) = \sum_{\omega \in \{0, \pm 1\}^d} (a_\omega \cos\langle \omega, x \rangle + b_\omega \sin\langle \omega, x \rangle) \quad (\text{C1})$$

of degree at most 1 over each one of d variables¹. Then for any point $\tilde{x} \in \mathbb{R}^d$, f can be written as

$$f(x) = f(\tilde{x}) + \sum_{i=1}^d T_i(x, \tilde{x}) + \varepsilon, \quad (\text{C2})$$

where T_i is defined as

$$T_i(x, y) := \sin(x_i - y_i) \left. \frac{\partial f(x)}{\partial x_i} \right|_{x=y} - (1 - \cos(x_i - y_i)) \left. \frac{\partial^2 f(x)}{\partial x_i^2} \right|_{x=y}, \quad (\text{C3})$$

¹ Here we take multivariate polynomials to also allow for negative frequencies up to minus the degree.

and ε contains all the cross-derivative contributions, that is, terms of the form

$$\left(\prod_{i=1}^d \frac{\partial^{n_i}}{\partial x_i^{n_i}} \right) f(x) \Big|_{x=y} \quad (\text{C4})$$

with $|n|_0 \geq 2$. Here $n = (n_i)_{i=1}^d$ is the vector of derivative orders, and $|n|_0 \geq 2$ means that it has at least two different non-0 entries.

Proof. We show the identity directly by using the general formula for multivariate Taylor expansion

$$f(x) = \sum_{|n|_1=0}^{\infty} \frac{\partial^{|n|_1} f(x)}{\partial x_1^{n_1} \dots \partial x_d^{n_d}} \Big|_{x=\tilde{x}} \prod_{i=1}^d \frac{(x_i - \tilde{x}_i)^{n_i}}{n_i!} \quad (\text{C5})$$

around an arbitrary point $\tilde{x} \in \mathbb{R}^d$. Again, we call $n = (n_i)_i \in \mathbb{N}^d$ the vector of derivative orders, and $|n|_1 = \sum_{i=1}^d n_i$ the total derivative order across all variables. We use the known identities for even and odd order derivatives of trigonometric polynomials, to get

$$\frac{\partial^{2k} f}{\partial x_i^{2k}} = (-1)^{k+1} \frac{\partial^2 f}{\partial x_i^2}, \text{ for all } k \in \mathbb{N} \setminus \{0\}, \quad (\text{C6})$$

$$\frac{\partial^{2k+1} f}{\partial x_i^{2k+1}} = (-1)^k \frac{\partial f}{\partial x_i}, \text{ for all } k \in \mathbb{N}. \quad (\text{C7})$$

For these identities, we need only use the fact that $\omega_i^3 = \omega_i$ for each component ω_i of each frequency vector $\omega \in \{0, \pm 1\}^d$. Now we re-organize the terms of the Taylor formula as

$$f(x) = \left[\prod_{i=1}^d \left(\sum_{n_i=0}^{\infty} \frac{(x_i - \tilde{x}_i)^{n_i}}{n_i!} \frac{\partial^{n_i}}{\partial x_i^{n_i}} \right) \right] f(x) \Big|_{x=\tilde{x}}. \quad (\text{C8})$$

We split each of the sums of operators into three separate terms

$$\sum_{n_i=0}^{\infty} \frac{(x_i - \tilde{x}_i)^{n_i}}{n_i!} \frac{\partial^{n_i}}{\partial x_i^{n_i}} = \mathbb{I} + \sum_{n_i=1}^{\infty} \frac{(x_i - \tilde{x}_i)^{2n_i}}{(2n_i)!} \frac{\partial^{2n_i}}{\partial x_i^{2n_i}} + \sum_{n_i=0}^{\infty} \frac{(x_i - \tilde{x}_i)^{2n_i+1}}{(2n_i+1)!} \frac{\partial^{2n_i+1}}{\partial x_i^{2n_i+1}}. \quad (\text{C9})$$

With this grouping of terms, we focus on the terms which overall contain derivatives (of any order) with respect to a single variable, and call all other terms E , to obtain

$$\begin{aligned} f(x) &= \left[\prod_{i=1}^d \left(\mathbb{I} + \sum_{n_i=1}^{\infty} \frac{(x_i - \tilde{x}_i)^{2n_i}}{(2n_i)!} \frac{\partial^{2n_i}}{\partial x_i^{2n_i}} + \sum_{n_i=0}^{\infty} \frac{(x_i - \tilde{x}_i)^{2n_i+1}}{(2n_i+1)!} \frac{\partial^{2n_i+1}}{\partial x_i^{2n_i+1}} \right) \right] f(x) \Big|_{x=\tilde{x}} \quad (\text{C10}) \\ &= \left[\mathbb{I} + \sum_{i=1}^d \left(\sum_{n_i=1}^{\infty} \frac{(x_i - \tilde{x}_i)^{2n_i}}{(2n_i)!} \frac{\partial^{2n_i}}{\partial x_i^{2n_i}} + \sum_{n_i=0}^{\infty} \frac{(x_i - \tilde{x}_i)^{2n_i+1}}{(2n_i+1)!} \frac{\partial^{2n_i+1}}{\partial x_i^{2n_i+1}} \right) + E \right] f(x) \Big|_{x=\tilde{x}} \\ &= f(\tilde{x}) + \sum_{i=1}^d \left(\sum_{n_i=1}^{\infty} \frac{(x_i - \tilde{x}_i)^{2n_i}}{(2n_i)!} \frac{\partial^{2n_i} f(x)}{\partial x_i^{2n_i}} \Big|_{x=\tilde{x}} + \sum_{n_i=0}^{\infty} \frac{(x_i - \tilde{x}_i)^{2n_i+1}}{(2n_i+1)!} \frac{\partial^{2n_i+1} f(x)}{\partial x_i^{2n_i+1}} \Big|_{x=\tilde{x}} \right) + E(f(x)) \Big|_{x=\tilde{x}}. \end{aligned}$$

At this point, we only need to use the identities we introduced above, and we call $E(f(x)) \Big|_{x=\tilde{x}} = \varepsilon$,

$$\begin{aligned} f(x) &= f(\tilde{x}) + \sum_{i=1}^d \left((1 - \cos(x_i - \tilde{x}_i)) \frac{\partial^2 f(x)}{\partial x_i^2} \Big|_{x=\tilde{x}} + \sin(x_i - \tilde{x}_i) \frac{\partial f(x)}{\partial x_i} \Big|_{x=\tilde{x}} \right) + \varepsilon \quad (\text{C11}) \\ &= f(\tilde{x}) + \sum_{i=1}^d T_i(x, \tilde{x}) + \varepsilon, \end{aligned}$$

which is what we had to show. \square

Appendix D: Step-by-step derivation of a digital twin neural network

We consider an encoding-first parametrized quantum circuit as a three step process:

1. Create a data-dependent quantum state $x \mapsto \rho(x)$.
2. Create a task-dependent parametrized observable $\vartheta \mapsto \mathcal{M}(\vartheta)$.
3. Take their inner product $f(x, \vartheta) = \text{tr}\{\rho(x)\mathcal{M}(\vartheta)\}$.

In this section, we derive a neuralization of such a PQC, in what we have dubbed a digital twiNN (twin Neural Network). Our approach corresponds to rewriting the PQC as a Neural Network, which involves rethinking the quantum objects $\rho(x)$ and $\mathcal{M}(\vartheta)$ as real-valued, large matrices. Again, we consider a three step process:

1. Create a data-dependent feature matrix $x \mapsto A(x)$.
2. Create a task-dependent linear layer $\vartheta \mapsto M(\vartheta)$.
3. Take their inner product $f(x, \vartheta) = \frac{1}{2} \text{tr}\{A(x)M(\vartheta)\}$.

We first introduce the recipe to construct A and M , and then show that the PQC and the twiNN indeed produce the same input-parameters-output relations. Keeping implementation in mind, we have different needs that A and M must fulfill. For A , we would like to have a closed-form expression that depends explicitly on x for each entry of the feature map. For M , we are satisfied with a formula that depends explicitly on each entry of $\mathcal{M}(\vartheta)$, thus being satisfied with only implicit dependence on the entries of ϑ .

We first introduce a natural way of expanding complex-valued matrices into slightly larger real-valued matrices in a way that respects the inner products. For an N -dimensional matrix $\mathbb{C}^{N \times N}$, consider the map

$$\mathbb{M}: \mathbb{C}^{N \times N} \rightarrow \mathbb{R}^{2N \times 2N}, \quad (\text{D1})$$

$$U \mapsto \begin{pmatrix} \Re(U) & -\Im(U) \\ \Im(U) & \Re(U) \end{pmatrix}, \quad (\text{D2})$$

where \Re and \Im stand for the entrywise real and imaginary parts, respectively.

Lemma 2 (The map \mathbb{M} is an isomorphism). *For any $U, V \in \mathbb{C}^{N \times N}$ complex matrices and UV their product, it follows that $\mathbb{M}(U)\mathbb{M}(V) = \mathbb{M}(UV)$.*

Proof. We show the identity directly. First we identify the real and imaginary parts of UV based on those of U and V , and then we show that the block structure of $\mathbb{M}(UV)$ corresponds to the correct parts.

Consider the expansion of U and V in their real and imaginary parts $U = \Re(U) + i\Im(U)$, $V = \Re(V) + i\Im(V)$. Then take their product $UV = \Re(U)\Re(V) - \Im(U)\Im(V) + i(\Re(U)\Im(V) + \Im(U)\Re(V))$, from where it follows that the real and imaginary parts of the product are $\Re(UV) = \Re(U)\Re(V) - \Im(U)\Im(V)$ and $\Im(UV) = \Re(U)\Im(V) + \Im(U)\Re(V)$. Next, consider the matrix-matrix multiplication

$$\begin{aligned} \mathbb{M}(U)\mathbb{M}(V) &= \begin{pmatrix} \Re(U) & -\Im(U) \\ \Im(U) & \Re(U) \end{pmatrix} \begin{pmatrix} \Re(V) & -\Im(V) \\ \Im(V) & \Re(V) \end{pmatrix} \\ &= \begin{pmatrix} \Re(U)\Re(V) - \Im(U)\Im(V) & -\Re(U)\Im(V) - \Im(U)\Re(V) \\ \Re(U)\Im(V) + \Im(U)\Re(V) & \Im(U)\Im(V) + \Re(U)\Re(V) \end{pmatrix} \\ &= \begin{pmatrix} \Re(UV) & -\Im(UV) \\ \Im(UV) & \Re(UV) \end{pmatrix} = \mathbb{M}(UV). \end{aligned} \quad (\text{D3})$$

Indeed, we have shown that $\mathbb{M}(U)\mathbb{M}(V) = \mathbb{M}(UV)$. □

Lemma 3 (Trace of $\mathbb{M}(H)$). *For any $H \in \mathbb{C}^{N \times N}$ Hermitian, $H = H^\dagger$, it holds that $\text{tr}\{\mathbb{M}(H)\} = 2 \text{tr}\{H\}$.*

Proof. We show this quickly and directly. One basic property of Hermitian matrices is that their main diagonal is real valued. Since the trace only involves the entries on the main diagonal, it follows that $\text{tr}\{H\} = \text{tr}\{\Re(H)\}$. From this, if we look at the block structure of $\mathbb{M}(H)$, it follows that the main diagonal of $\mathbb{M}(H)$ is nothing but two copies of the main diagonal of H , one after the other. And, hence, it is clear that $\text{tr}\{\mathbb{M}(H)\} = 2 \text{tr}\{\Re(H)\} = 2 \text{tr}\{H\}$, thus proving the lemma. □

Next, for the feature map $A(x)$. Without loss of generality, we assume the PQC encodes d -dimensional data $x \in \mathbb{R}^d$ in d qubits, each component being introduced once as the parameter of a parametrized Pauli- X rotation $R_X(x_i)$, on the corresponding qubit.

Lemma 4 (Entries of A). *We have efficient formulas for computing any entry of $A(x)$.*

Proof. We give the formula in the following. The encoding gate $U(x)$ is of the form $U(x) = \bigotimes_{j=1}^d R_X(x_j)$ and, when applied on the $|0\rangle$ state vector, it produces a data-dependent state as

$$\begin{aligned}
U(x)|0\rangle\langle 0|U(x)^\dagger &= \bigotimes_{j=1}^d R_X(x_j)|0\rangle\langle 0|R_X(x_j)^\dagger \\
&= \bigotimes_{j=1}^d \begin{pmatrix} \cos \frac{x_j}{2} & \\ -i \sin \frac{x_j}{2} & \end{pmatrix} \begin{pmatrix} \cos \frac{x_j}{2} & i \sin \frac{x_j}{2} \\ & \end{pmatrix} \\
&= \bigotimes_{j=1}^d \begin{pmatrix} \cos^2 \frac{x_j}{2} & i \cos \frac{x_j}{2} \sin \frac{x_j}{2} \\ -i \cos \frac{x_j}{2} \sin \frac{x_j}{2} & \sin^2 \frac{x_j}{2} \end{pmatrix} \\
&= \bigotimes_{j=1}^d \begin{pmatrix} \frac{1+\cos x_j}{2} & \frac{i \sin x_j}{2} \\ -\frac{i \sin x_j}{2} & \frac{1-\cos x_j}{2} \end{pmatrix}.
\end{aligned} \tag{D4}$$

At this point, we introduce short-hand notation

$$c_j^+ := 1 + \cos x_j, \tag{D5}$$

$$c_j^- := 1 - \cos x_j, \tag{D6}$$

$$s_j := \sin x_j. \tag{D7}$$

We also introduce bit-string indices $k, l \in \{0, 1\}^d$, whose entries are indexed by j , $k = (k_j)_j$, with $k_j \in \{0, 1\}$, and similarly for l . We make use of the kronecker delta $\delta_{a,b}$, which is the indicator function for the condition $a = b$. Introducing these new symbols in the equations we obtain a formula for the entries of $\rho(x)$ which should be particularly friendly to implement on a classical computer,

$$\begin{aligned}
U(x)|0\rangle\langle 0|U(x)^\dagger &= \dots = \frac{1}{2^d} \bigotimes_{j=1}^d \begin{pmatrix} c_j^+ & i s_j \\ -i s_j & c_j^- \end{pmatrix} \\
&= \frac{1}{2^d} \bigotimes_{j=1}^d [c_j^+ |0\rangle\langle 0| + c_j^- |1\rangle\langle 1| + i s_j (|0\rangle\langle 1| + |1\rangle\langle 0|)] \\
&= \frac{1}{2^d} \bigotimes_{j=1}^d \left(\sum_{k_j, l_j=0}^1 [\delta_{k_j, l_j} (c_j^+ \delta_{k_j, 0} + c_j^- \delta_{k_j, 1}) + (1 - \delta_{k_j, l_j}) (-1)^{k_j} s_j] |k_j\rangle\langle l_j| \right) \\
&= \frac{1}{2^d} \sum_{k, l \in \{0, 1\}^d} \left(\prod_{j=1}^d \delta_{k_j, l_j} (c_j^+ \delta_{k_j, 0} + c_j^- \delta_{k_j, 1}) + (1 - \delta_{k_j, l_j}) (-1)^{k_j} s_j \right) |k\rangle\langle l| \\
&= \rho(x).
\end{aligned} \tag{D8}$$

In order to obtain $A(x) = M(\rho(x))$, we need to separate the real and imaginary parts of $\rho(x)$. Luckily, each entry of $\rho(x)$ is either pure real or pure imaginary, with no mixed complex entries. Upon visual inspection of the formula for the entries of $\rho(x)$, we observe a pattern that relates each of the trigonometric functions c_j^+ , c_j^- , s_j to each of the computational basis elements $|0/1\rangle\langle 0/1|$. We want to introduce a map which captures this pattern in the form of a look-up table

$$|0_j\rangle\langle 0_j| \mapsto c_j^+, \tag{D10}$$

$$|0_j\rangle\langle 1_j| \mapsto i s_j, \tag{D11}$$

$$|1_j\rangle\langle 0_j| \mapsto -i s_j, \tag{D12}$$

$$|1_j\rangle\langle 1_j| \mapsto c_j^-. \tag{D13}$$

We call it g_j for ease of notation, and its formula reads as

$$g_j: \{0, 1\}^2 \rightarrow \{c_j^+, c_j^-, s_j\}, \quad (\text{D14})$$

$$k_j, l_j \mapsto g_j(k_j, l_j) = i^{3k_j+l_j} \left(\delta_{k_j, l_j} + \cos(x_j - \frac{\pi}{2}(\delta_{k_j,1} + \delta_{l_j,1})) \right). \quad (\text{D15})$$

Similarly, we can introduce a global map g which corresponds to g_j on each qubit, so for any $k, l \in \{0, 1\}^d$,

$$\begin{aligned} g(k, l) &= \prod_{j=1}^d g_j(k_j, l_j) \\ &= \prod_{j=1}^d i^{3k_j+l_j} \left(\delta_{k_j, l_j} + \cos(x_j - \frac{\pi}{2}(\delta_{k_j,1} + \delta_{l_j,1})) \right) \\ &= i^{\sum_{j=1}^d 3k_j+l_j} \prod_{j=1}^d \left(\delta_{k_j, l_j} + \cos(x_j - \frac{\pi}{2}(\delta_{k_j,1} + \delta_{l_j,1})) \right) \\ &= i^{3|k|+|l|} \prod_{j=1}^d \left(\delta_{k_j, l_j} + \cos(x_j - \frac{\pi}{2}(\delta_{k_j,1} + \delta_{l_j,1})) \right), \end{aligned} \quad (\text{D16})$$

We used the notation $|k|$ for the parity of the bit-string k (the number of 1s). With this we reached yet another implementation-friendly formula for the entries of $\rho(x)$, and in fact one where we can very quickly extract the real and imaginary parts

$$\rho(x) = \sum_{k, l \in \{0, 1\}^d} \frac{g(k, l)}{2^d} |k\rangle\langle l|, \quad (\text{D17})$$

$$\begin{aligned} \Re(\rho(x)) &= \frac{1}{2^d} \sum_{k, l \in \{0, 1\}^d} \Re(g(k, l)) |k\rangle\langle l| \\ &= \frac{1}{2^d} \sum_{k, l \in \{0, 1\}^d} \delta_{|k \oplus l, 0} g(k, l) |k\rangle\langle l|, \end{aligned} \quad (\text{D18})$$

$$\begin{aligned} \Im(\rho(x)) &= \frac{1}{2^d} \sum_{k, l \in \{0, 1\}^d} \Im(g(k, l)) |k\rangle\langle l| \\ &= \frac{1}{2^d} \sum_{k, l \in \{0, 1\}^d} \delta_{|k \oplus l, 1} (-i) g(k, l) |k\rangle\langle l|. \end{aligned} \quad (\text{D19})$$

What these formulas relate is that the ‘‘global even’’ terms (taking addition modulo 2 for both bit-strings $k \oplus l$) are pure real, and the ‘‘global odd’’ terms are pure imaginary.

Bringing everything together, we are left with a formula for computing each entry of the data-dependent matrix $A(x)$, that corresponds to the first layer of our twiNN. If we call $G(x)$ the matrix whose entries are $(g(k, l))_{k, l \in \{0, 1\}^d}$, we obtain

$$A(x) = \frac{1}{4^n} \begin{pmatrix} \Re(G(x)) & -\Im(G(x)) \\ \Im(G(x)) & \Re(G(x)) \end{pmatrix}. \quad (\text{D20})$$

For any entry of $A(x)$, we need only compute $g(k, l)$ for the corresponding bit-strings $k, l \in \{0, 1\}^d$, whose complexity is linear in d , so the formula is efficient for any entry. \square

Finally, for the task-dependent linear layer, we take $M(\vartheta) = \mathbf{M}(\mathcal{M}(\vartheta))$. In order to compute the entries of $\mathcal{M}(\vartheta)$, we need to have a decomposition of $\mathcal{M}(\vartheta)$ as a fixed observable and a variational circuit $\mathcal{M}(\vartheta) = V(\vartheta)^\dagger \mathcal{M}_0 V(\vartheta)$. This makes sense for example if we assume the fixed Hamiltonian admits an efficient classical description, and if the variational block is composed of several parametrized local unitaries, for which we also have efficient classical descriptions.

Lemma 5 (Expectation values). *With the above definitions, it holds that*

$$\text{tr}\{\rho(x)\mathcal{M}(\vartheta)\} = \frac{1}{2} \text{tr}\{A(x)M(\vartheta)\}. \quad (\text{D21})$$

Proof. We prove this statement directly. Since we have $A(x) = \mathbf{M}(\rho(x))$ and $M(\vartheta) = \mathbf{M}(\mathcal{M}(\vartheta))$, it follows from Lemma 2 that $A(x)M(\vartheta) = \mathbf{M}(\rho(x)\mathcal{M}(\vartheta))$. Subsequently, Lemma 3 confirms that $\text{tr}\{A(x)M(\vartheta)\} = 2 \text{tr}\{\rho(x)\mathcal{M}(\vartheta)\}$. \square

Appendix E: Derivation of quantum layerwise relevance propagation for the twiNN, linear rule, and encoding rule

The twiNN produces an output as a two-step process

$$x \mapsto A \mapsto f = \text{tr}\{AM\}. \quad (\text{E1})$$

We drop the explicit dependence on x and ϑ for ease of notation. Correspondingly, we are after a relevance attribution method that works in two steps, in reversed order. For a given tensor (real number, vector, or matrix) T , we use the notation $R(T)$ for the relevance attribution of T . The relevance tensor $R(T)$ has the same shape and size as T , and we denote its components with the same set of indices. For example, for a matrix with two indices T_{ab} , the relevance of the (a, b) entry of T is denoted as $R_{ab}(T)$.

With this, we propose a relevance propagation algorithm that, at every step, takes into account the relevance tensor of the previous step, and the tensors of that and the previous step. Diagrammatically, this approach looks like follows:

$$\begin{array}{ccccc}
 R(f) & \xrightarrow{\text{linear rule}} & R(A) & \xrightarrow{\text{encoding rule}} & R(x) \\
 \uparrow & & \uparrow & \nearrow & \uparrow \\
 f & \xleftarrow{\text{linear step}} & A & \xleftarrow{\text{encoding step}} & x
 \end{array}$$

In this section we derive the formula for the linear rule and for the encoding rule. As is usual, we start from $R(f) = f$.

1. Linear rule

The linear rule is almost immediate. We just need to look at the formula for f in terms of A , and then visual inspection gives us the rule

$$f(A) = \text{tr}\{AM\} = \sum_{i,j} A_{i,j} M_{i,j}. \quad (\text{E2})$$

Since f is a linear function of the entries of A , the relevance associated to the entries of A is immediately read out: $R_{i,j}(A) = A_{i,j} M_{i,j}$. The linear step ensures exact conservation by construction $R(f) = \sum_{i,j} R_{i,j}(A)$.

2. Encoding rule

The encoding step is non-linear, so we have a bit of work to do. From the formula for the entries of A in terms of x , we realize that each entry of A is almost a trigonometric monomial of degree 1 on each of the components of x . Only the 1 summands contained in the c^\pm terms cause the product to turn into a sum of several monomial terms. Yet, for the function basis $t_k \in \{c_j^\pm, s_j\}$, we have that each entry $A_{i,j}$ is either 0 or of the form

$$A_{i,j} = \frac{1}{2^d} \prod_{k=1}^d t_k^{(i,j)}. \quad (\text{E3})$$

The crucial fact here is that each $A_{i,j}$ entry as a function is a trigonometric polynomial of degree at most 1 on each of the d components of x . Then, following Proposition 1, we have seen they admit the following approximation, for any *root point* $\tilde{x}^{(i,j)}$,

$$A_{i,j}(x) \approx A_{i,j}(\tilde{x}^{(i,j)}) + \sum_{k=1}^d T_k^{(i,j)}(x, \tilde{x}^{(i,j)}) + \varepsilon, \quad (\text{E4})$$

with

$$T_k^{(i,j)}(x, \tilde{x}^{(i,j)}) = \sin(x_k - \tilde{x}_k^{(i,j)}) \left. \frac{\partial A_{i,j}(x)}{\partial x_k} \right|_{x=\tilde{x}^{(i,j)}} + (1 - \cos(x_k - \tilde{x}_k^{(i,j)})) \left. \frac{\partial^2 A_{i,j}(x)}{\partial x_k^2} \right|_{x=\tilde{x}^{(i,j)}}, \quad (\text{E5})$$

and where ε contains all cross-order derivative terms. With this approximation, and assuming we have suitable root points for each (i, j) , for which $A_{i,j}(\tilde{x}^{(i,j)}) \approx 0$, then the encoding rule reads

$$R_k(x) = \sum_{i,j} T_k^{(i,j)}(x, \tilde{x}^{(i,j)}) \frac{R_{i,j}(A)}{A_{i,j}}. \quad (\text{E6})$$

This rule is approximately conservative by construction, with the amount of conservation being based on the quality of approximation of the Fourier decomposition

$$\begin{aligned}
\sum_k R_k(x) &= \sum_k \sum_{i,j} T_k^{(i,j)}(x, \tilde{x}^{(i,j)}) \frac{R_{i,j}(A)}{A_{i,j}} \\
&= \sum_{i,j} \left(\sum_k T_k^{(i,j)}(x, \tilde{x}^{(i,j)}) \right) \frac{R_{i,j}(A)}{A_{i,j}} \\
&\approx \sum_{i,j} A_{i,j} \frac{R_{i,j}(A)}{A_{i,j}} \\
&= \sum_{i,j} R_{i,j}(A).
\end{aligned} \tag{E7}$$

In order to implement the encoding rule, one needs to evaluate not only each entry of $A(x)$, but also the first and second partial derivatives of $A(x)$ with respect to each component of x . Fortunately, these are really straightforward to evaluate, again due to the simple trigonometric structure of $A(x)$ we have been exploiting. In particular, we need only use the parameter-shift rule. We give it first for an individual qubit, and then see how it generalizes right away to several qubits, provided each qubit uploads a different component

$$\rho^1(x_1) = \frac{1}{2} \begin{pmatrix} 1 + \cos(x_1) & i \sin(x_1) \\ -i \sin(x_1) & 1 - \cos(x_1) \end{pmatrix}, \tag{E8}$$

$$\begin{aligned}
\frac{d\rho^1(x_1)}{dx_1} &= \frac{1}{2} \begin{pmatrix} -\sin(x_1) & i \cos(x_1) \\ -i \cos(x_1) & +\sin(x_1) \end{pmatrix} \\
&= \frac{\rho^1(x_1 + \frac{\pi}{2}) - \rho^1(x_1 - \frac{\pi}{2})}{2}.
\end{aligned} \tag{E9}$$

Here $\rho^1(x_1)$ represented a single-qubit encoded state. The case is not much harder for the whole d -qubit state

$$\rho(x) = \rho^1(x_1) \otimes \cdots \otimes \rho^d(x_d), \tag{E10}$$

$$\begin{aligned}
\partial_k [\rho(x)] &= \partial_k [\rho^1(x_1) \otimes \cdots \otimes \rho^k(x_k) \otimes \cdots \otimes \rho^d(x_d)] \\
&= \rho^1(x_1) \otimes \cdots \otimes \partial_k \rho^k(x_k) \otimes \cdots \otimes \rho^d(x_d) \\
&= \rho^1(x_1) \otimes \cdots \otimes \frac{\rho^k(x_k + \frac{\pi}{2}) - \rho^k(x_k - \frac{\pi}{2})}{2} \otimes \cdots \otimes \rho^d(x_d) \\
&= \frac{\rho(x + \frac{\pi}{2} \hat{e}_k) - \rho(x - \frac{\pi}{2} \hat{e}_k)}{2},
\end{aligned} \tag{E11}$$

where \hat{e}_k denotes the k^{th} basis vector. The formula can be used recursively for higher-order derivatives

$$\begin{aligned}
\partial_k^2 \rho(x) &= \partial_k [\partial_k \rho(x)] \\
&= \partial_k \left[\frac{\rho(x + \frac{\pi}{2} \hat{e}_k) - \rho(x - \frac{\pi}{2} \hat{e}_k)}{2} \right] \\
&= \partial_k \left[\frac{\rho(x + \pi \hat{e}_k) - \rho(x) - \rho(x) + \rho(x - \pi \hat{e}_k)}{4} \right] \\
&= -\frac{\rho(x) - \rho(x + \pi \hat{e}_k)}{2},
\end{aligned} \tag{E12}$$

in the last line we used the 2π -periodicity of $\rho(x)$ to say $\rho(x + \pi \hat{e}_k) = \rho(x - \pi \hat{e}_k)$. The identities we found for $\rho(x)$ apply directly to $A(x)$, since the map from $\rho(x)$ to $A(x)$ is linear. That means also we do not need to give explicit closed forms for the derivatives of the $g(k, l)$ functions introduced above, since we have shown in order to compute those we need only evaluate the functions themselves on shifted locations.

So, there is only one piece missing in this puzzle, and that is how we find the root points $\tilde{x}^{(i,j)}$. Note that the root points in this case are slightly different in nature as the root points we have talked about in other approaches. Here we are not looking for a single point where the overall function f nullifies. Rather, we are looking for as many points as components A has (exponential

in the input dimension). Luckily, since we have closed formulas for the entries of A , we can take advantage of the shared structure among all of them.

Each of the entries is a product of trigonometric functions of a single variable

$$A_{i,j}(x) = \prod_k t_k^{(i,j)}(x_k). \quad (\text{E13})$$

So, in order to achieve $A_{i,j}(\tilde{x}^{(i,j)}) = 0$, it suffices to find a zero of the trigonometric function $t_k^{(i,j)}$ corresponding to a single component $t_k^{(i,j)}(\tilde{x}_k^{(i,j)}) = 0$.

With this knowledge, we here propose an algorithm to find root points for each component. Because the property that must be fulfilled by the root point only affects a single component, it could be that a single point is a good root point for several entries of A simultaneously. In the algorithm we propose we concentrate on finding the *closest* root point to a given x for each $A_{i,j}$ in the Euclidean sense. We first give the procedure in plain language, and at the end we attach the pseudo-code. As we did in the proof of Lemma 4, we consider again the map

$$|0_k\rangle\langle 0_k| \mapsto c_k^+, \quad (\text{E14})$$

$$|0_k\rangle\langle 1_k| \mapsto i s_k, \quad (\text{E15})$$

$$|1_k\rangle\langle 0_k| \mapsto -i s_k, \quad (\text{E16})$$

$$|1_k\rangle\langle 1_k| \mapsto c_k^-. \quad (\text{E17})$$

From here, we read the corresponding root points for each trigonometric function

$$|0_k\rangle\langle 0_k| \mapsto 1 + \cos(\tilde{x}_k) = 0 \iff \tilde{x}_k = \pm\pi, \quad (\text{E18})$$

$$|0_k\rangle\langle 1_k| \mapsto i \sin(\tilde{x}_k) = 0 \iff \tilde{x}_k \in \{0, \pm\pi\}, \quad (\text{E19})$$

$$|1_k\rangle\langle 0_k| \mapsto -i \sin(\tilde{x}_k) = 0 \iff \tilde{x}_k \in \{0, \pm\pi\}, \quad (\text{E20})$$

$$|1_k\rangle\langle 1_k| \mapsto 1 - \cos(\tilde{x}_k) = 0 \iff \tilde{x}_k = 0. \quad (\text{E21})$$

Each of the conditions defines a high-dimensional grid (a partition via periodical hyperplanes in every direction) on input space. The question becomes, given x and (i, j) , what is the closest root point $\tilde{x}^{(i,j)}$?

From x , we start by building three related vectors $x^{(1)} = x$, $x^{(2)} = x + \pi\vec{1}$, and $x^{(3)} = x - \pi\vec{1}$. Here $\vec{1}$ represents a vector of ones. Now, we consider the matrix $(x_m^{(n)})_{\substack{m \in \{1, \dots, d\} \\ n=1,2,3}}$, and consider the sequence $(m_l, n_l)_{l \in \{1, \dots, 3d\}}$, which corresponds to the sorting of the elements of $(x_m^{(n)})_{\substack{m \in \{1, \dots, d\} \\ n=1,2,3}}$, according to their magnitude. For instance, the first few elements are

$$(m_1, n_1) = \arg \min_{(m,n) \in \{1, \dots, d\} \times \{1,2,3\}} \{|x_m^{(n)}|\}, \quad (\text{E22})$$

$$(m_2, n_2) = \arg \min_{(m,n) \in \{1, \dots, d\} \times \{1,2,3\} \setminus \{(m_1, n_1)\}} \{|x_m^{(n)}|\}, \quad (\text{E23})$$

⋮

$$(m_l, n_l) = \arg \min_{(m,n) \in \{1, \dots, d\} \times \{1,2,3\} \setminus \bigcup_{s=1}^{l-1} \{(m_s, n_s)\}} \{|x_m^{(n)}|\}, \quad (\text{E24})$$

⋮

Now, take (m_1, n_1) , which identifies the hyperplane to which x is closest. The first index m_1 refers to the component x_{m_1} , and the second index n_1 identifies whether x_{m_1} is close to 0, π , or $-\pi$. We identify the specific root point $\tilde{x}^{(i,j)} = (\tilde{x}_k^{(i,j)})_k$ from the value of n_1 :

$$\tilde{x}_k = \begin{cases} x_k & \text{for } k \neq m_1, \\ \begin{cases} 0 & \text{if } n_1 = 1, \\ -\pi & \text{if } n_1 = 2, \\ \pi & \text{if } n_1 = 3, \end{cases} & \text{for } k = m_1. \end{cases} \quad (\text{E25})$$

Indeed, x differs from \tilde{x} only on a single component, $k = m_1$. We have identified the root point, but we have not said which component this is a root point of. For that, we go back to the identification above, and again reach an answer based on the value

of n_1 . Instead of starting from an (i, j) and finding the matching root point, we go the other way: we start from the root point we just defined, and check for which choices of (i, j) it is a valid root point.

We consider the bit-string representation of i and j , and then focus on the m_1^{th} bits i_{m_1}, j_{m_1} . We reverse the correspondence from above, now taking the indexing of n_1 , for which we have the look-up table

$$n_1 \mapsto |i_{m_1}\rangle\langle j_{m_1}|, \quad (\text{E26})$$

$$1 \mapsto \{|0\rangle\langle 1|, |1\rangle\langle 0|, |1\rangle\langle 1|\}, \quad (\text{E27})$$

$$\{2, 3\} \mapsto |0\rangle\langle 0|. \quad (\text{E28})$$

The way to read this map is: if $n_1 = 1$, then \tilde{x} as defined above is the root point for all choices of (i, j) where the m_1^{th} bit of either i or j is 1. Conversely, if n_1 is either 2 or 3, then \tilde{x} as defined above is the root point for all choices of (i, j) where the m_1^{th} bit of both i and j is 0.

After this step, we have found a root point that works for either 1/4 or 3/4 of all choices of (i, j) . Notice in the earlier paragraph we talk about *all* choices of (i, j) for which a particular bit is a particular value.

After this step, we proceed to use the next element in the sequence, (m_2, n_2) . If it happens that $m_2 = m_1$, then there is a chance this is the last step. If either n_1 or n_2 are 1, then it follows that \tilde{x} as defined above is the root point for all choices of (i, j) . If neither of them is 1, then the point corresponding to (m_2, n_2) does not provide new root points, and we move further up the sequence.

If, on the contrary, we have that $m_2 \neq m_1$, then we find a new root point. In this case, we repeat all the steps replacing (m_1, n_1) by (m_2, n_2) up until the point where we identify for which choices of (i, j) the newly defined \tilde{x} is a valid root point. As by now we already had a fraction of all possible choices associated with the first root point, now we only look at the remaining ones, those for which we do not have a root point yet. At this point, in order to see which choices of (i, j) have the new \tilde{x} as a root point, we will check the m_2^{th} bit of i and j . Again, this results on finding a root point for either one or three quarters of the remaining choices of (i, j) .

As one can expect, we need to repeat this process at most $\mathcal{O}(d)$ times before we are guaranteed to have found a root point for each possible choice of (i, j) . The search finishes only when a new m_l is equal to one of the previous ones, with either $n_l = 1$ or the previous one being equal to 1.

Algorithm 1 Root points for the encoding rule**Require:**1: $A \in \mathbb{R}^{2^d \times 2^d}$

▷ Real-valued matrix for data-dependent state.

2: $x \in \mathbb{R}^d$

▷ Input vector.

Ensure: $\tilde{x} \in \mathbb{R}^d \times \mathbb{R}^{2^d \times 2^d}$ vector of root points for each entry of A .

3:

4: **for** $i \in \{0, 1\}^d, j \in \{0, 1\}^d$ **do**5: $\tilde{x}_k^{(i,j)} \leftarrow \text{void}$

▷ Initialize empty root point vector.

6: **end for**

7:

8: $x^{(1)} \leftarrow x$ 9: $x^{(2)} \leftarrow x + \pi \vec{1}$ 10: $x^{(3)} \leftarrow x - \pi \vec{1}$ 11: $(M, N) \leftarrow \text{sort}(m, n)$ according to ascending $|x_m^{(n)}|$

▷ Sort indices according to distance to grid.

12: **for** $(m_l, n_l) \in (M, N)$ **do**13: **if** $n_l = 1$ **then**14: **for** $(i, j) \in \{(i_{m_l}, j_{m_l}) \in \{(0, 1), (1, 0), (1, 1)\}\}$ **do**

▷ Check 3/4 of all remaining bitstrings.

15: **if** $\tilde{x}^{(i,j)} = \text{void}$ **then**

▷ Assign root point only if it had not already been assigned.

16: $\tilde{x}_k^{(i,j)} \leftarrow \begin{cases} 0 & \text{if } k = m_l \\ x_k & \text{else} \end{cases}$ ▷ Assign root point differing from x only on one component.17: **end if**18: **end for**19: **else** $n_l \in \{2, 3\}$ 20: **for** $(i, j) \in \{(i_{m_l}, j_{m_l}) = (0, 0)\}$ **do**

▷ Check 1/4 of all remaining bitstrings.

21: **if** $\tilde{x}^{(i,j)} = \text{void}$ **then**

▷ Assign root point only if it had not already been assigned.

22: $\tilde{x}_k^{(i,j)} \leftarrow \begin{cases} -\pi & \text{if } n_l = 2 \\ \pi & \text{if } n_l = 3 \\ x_k & \text{else} \end{cases}$ if $k = m_l$ ▷ Assign root point differing from x only on one component, depending on n_l .23: **end if**24: **end for**25: **end if**26: **for** $i \in \{0, 1\}^d, j \in \{0, 1\}^d$ **do**

▷ Early-stopping criterion.

27: **if** $\tilde{x}^{(i,j)} \neq \text{void}$ **then**▷ Check if all entries of A have a root point assigned to them at the end of each iteration.28: **return** \tilde{x}

▷ If all entries have a root point, the algorithm is finished.

29: **end if**30: **end for**31: **end for**32: **return** \tilde{x}

As a potential limitation, this algorithm finds root points which differ from the original point in only one direction. That means we are restricting the relevance propagation rule to distribute relevance only onto one component. This is a negative property because the component onto which relevance gets distributed is not decided based on task-dependent information, but rather only on the geometry of the encoding functions that arise from the encoding step.

Appendix F: Evaluation metrics

The first evaluation metric measures the alignment between the proposed explanation and the mask, and we call it *explanation alignment* Q_A . To have a well-behaved score, we compute the fraction of the absolute relevance that is assigned to the correct components:

$$Q_A(x) = \frac{\sum_{i=1}^d |E_i(x)| M_i(x)}{\sum_{i=1}^d |E_i(x)|}. \quad (\text{F1})$$

We note that whether $M_i(x)$ is 1 or 0 depends on the correct label for x . Also, the explanation alignment is indirectly related to the conservativity of the explanation $E(x)$, as even if $\sum_i E_i(x) \approx f(x)$, here we sum over the absolute values of the explanation.

The second quality-of-explanation metric considers a different kind of alignment between the explanation E and the mask M , namely their correlation coefficient, to reach the *Pearson correlation* metric Q_P . This metric deals not in absolute values, but in

deviations from the mean, also normalized to produce a well-behaved score:

$$Q_P(x) = \text{Corr}_{i \in [d]}(E_i(x), M_i(x)) = \frac{\mathbb{E}_i [(E_i(x) - \mathbb{E}_j [E_j(x)])(M_i(x) - \mathbb{E}_j [M_j(x)])]}{\sqrt{\mathbb{E}_i [E_i(x) - \mathbb{E}_j [E_j(x)]]^2 \mathbb{E}_i [M_i(x) - \mathbb{E}_j [M_j(x)]]^2}}, \quad (\text{F2})$$

here we took $[d] = \{1, \dots, d\}$, and \mathbb{E}_i means the expectation value over the indices $i \in [d]$.

Finally, we introduce a quality-of-explanation metric using the framework of *receiver operator characteristics* (ROC). This metric produces a score given a set of explanations, and not for a single input. Introducing a real-valued threshold $\alpha \in [0, 1]$, we momentarily define as “well-explained” all inputs whose explanation alignment is higher than α , and as “wrong-explained” all inputs whose alignment to the unimportant components (for example the opposite of the mas $1 - M(x)$) is higher than α . This way, an input could be well-explained and wrong-explained at the same time, for the purposes of this metric.

Computing this third-and-last evaluation metric follows a two-step process. First, given a set of inputs S we compute the fraction of well-explained $r_+(\alpha)$ and wrong-explained $r_-(\alpha)$ inputs, for different values $\alpha \in [0, 1]$:

$$r_+(\alpha) = \frac{|\{x \in S \mid x \text{ well-explained}\}|}{|S|} \quad (\text{F3})$$

$$r_-(\alpha) = \frac{|\{x \in S \mid x \text{ wrong-explained}\}|}{|S|}, \quad (\text{F4})$$

and from here, we consider the parametrized curve $\gamma(\alpha) = (r_-(\alpha), r_+(\alpha))$ for $\alpha \in [0, 1]$. From this parametrized curve, we numerically extract the *receiver operator characteristic* curve (ROC), where the dependence between r_+ and r_- is made explicit, removing the parameter α . That is, we take the fraction of well-explained inputs as a function of the fraction of wrong-explained inputs $r_+(r_-) \in [0, 1]$, for $r_- \in [0, 1]$. Next, we compute the *area under the curve* (AUC) for the ROC curve, which gives us the metric:

$$Q_{\text{ROC}}(S) = \int_0^1 r_+(r_-) dr_-. \quad (\text{F5})$$

Given that both $r_+, r_- \in [0, 1]$, it follows that $Q_{\text{ROC}} \in [0, 1]$. The ROC curve has an interpretation coming from binary classification tasks, below we show an example of how such curves look like in practice. In the main text we keep only the area under the curve, and not the whole curve.

Comparing the ROC curves for different attribution methods thus allows for a qualitative assessment of which methods perform better at accurately attributing relevance to the main dimensions while minimizing irrelevant attributions. This behavior is then succinctly summarized in the area under the ROC curve (AUC) as a single number.

Appendix G: Experimental setup

In this part of the appendix, we give an overview of the employed data-set and task, the different XAI and XQML methods we compare, and the evaluation metrics we use.

1. Dataset and learning task

To compare the individual methods, we train a Parametrized Quantum Circuit (PQC) on the following setup. We generate a synthetic data-set for a 6-dimensional classification task with four classes. For each class, we select three main dimensions as follows:

Class	Main Dimensions	Remaining Dimensions
0	{0, 1, 2}	{3, 4, 5}
1	{3, 4, 5}	{0, 1, 2}
2	{0, 2, 4}	{1, 3, 5}
3	{1, 3, 5}	{0, 2, 4}

Table IV. Class specifications for main and remaining dimensions.

Next, we generate data for each class by sampling entries from one of two distributions, depending on the main dimensions. For each class, the main dimensions follow a Gaussian distribution, with different means for different classes. The remaining dimensions are sampled uniformly at random in an interval centered around the origin. Specifically, the data-generating distribution looks as follows:

$$(\mathbf{x}^{(c)}, y^{(c)}) = \left(\left\{ (x_1^{(c)}, x_2^{(c)}, x_3^{(c)}) \sim \mathcal{N}(\mu, R^2) \right\} \cup \left\{ (x_4^{(c)}, x_5^{(c)}, x_6^{(c)}) \sim \text{Unif}([-m, m]) \right\}, y^{(c)} = i \right), \quad (\text{G1})$$

$$\text{where } c \in \{0, 1, 2, 3\}. \quad (\text{G2})$$

We use superscripts to refer to inputs in the c^{th} class. To tackle this problem in a supervised learning fashion, we are given a training set of labeled data: $S = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ sampled according to the distribution above.

Thus, each class' data forms a cluster characterized by a specific subset of three dimensions being drawn from a normal distribution, while the remaining dimensions are sampled uniformly at random. The normal distribution is characterized by its mean $\mu \in \mathbb{R}^3$ and the standard deviation $\sigma \in \mathbb{R}$. Table IV below lists the main dimensions and remaining dimensions used for each class:

Distribution	Parameter	Value
Normal	Mean (μ)	$(\frac{1}{2}, \frac{1}{2}, 0)$
Normal	Standard deviation (σ)	$\frac{1}{\sqrt{2}} \times 0.2$
Uniform	Interval endpoint (m)	$\{0.1, 0.5, \pi\}$

Table V. Parameters for the normal and uniform distributions.

Table V lists the parameters for the normal and uniform distributions used in the data generation process. In total, we draw $n = 1000$ samples for each respective class.

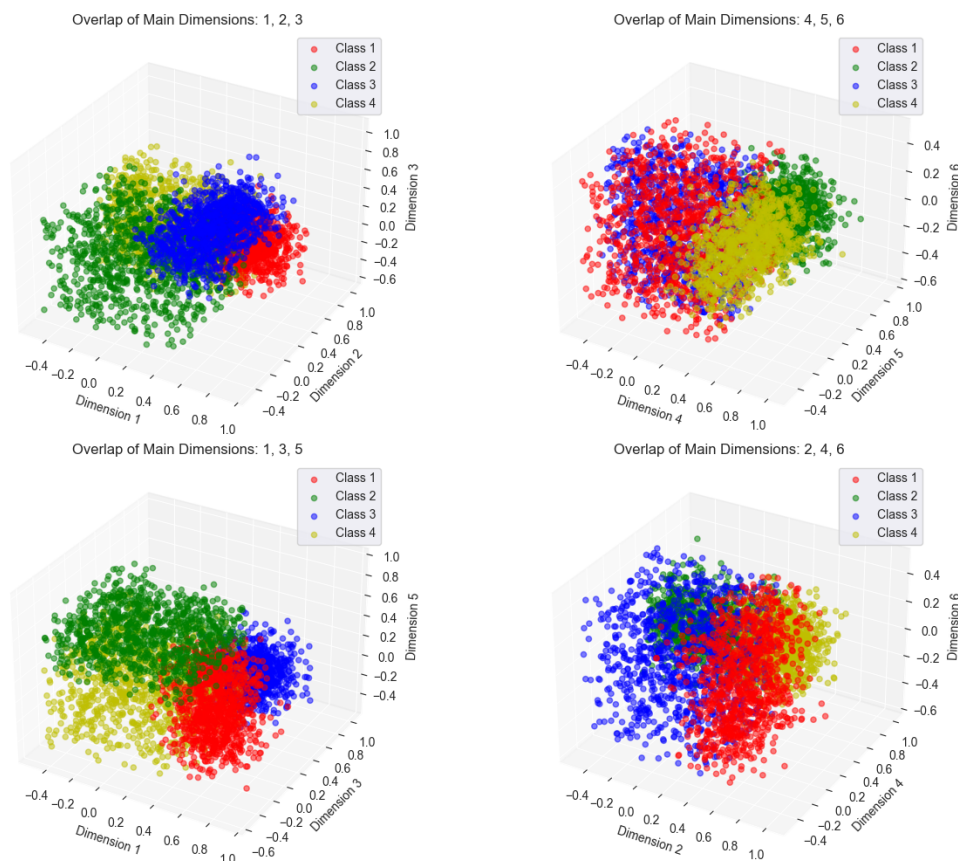


Figure 2. Visualization of the data-set across different subsets of \mathbb{R}^6 .

2. Parametrized quantum circuits

As an exemplary learning model, we employ a simple, rotational encoding based parameterized quantum circuit together with a trainable entangling gate set. The embedding of the numerical data is achieved by encoding the data points' individual values on differing qubits using Pauli- X rotations. The trainable part of the PQC is given by n_{layer} repetitions of blocks of strongly entangling layers. The latter consist of single qubit rotational gates followed by entanglers [93]. We compute the prediction of a label by measuring the expectation values of the single-qubit Pauli- Z observable on the first 4 qubits, respectively. An exemplary circuit for $n_{\text{layer}} = 5$ is shown in Fig. 3.

a. Training

We use a categorical cross-entropy loss function to train the model. For a given batch of input data and associated labels $\mathcal{B} = \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$ and resulting model output $\mathbf{z}_i = f(\mathbf{x}_i)$, it is given by

$$\mathcal{L}_{\text{CE}} = -\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C y_{ic} \log(\hat{y}_{ic} + \epsilon) \quad (\text{G3})$$

where N is the number of samples in the batch, C is the number of classes, $\epsilon = 10^{-10}$ is a small constant added for numerical stability, and

$$\hat{\mathbf{y}}_i = \text{softmax}(\mathbf{z}_i), \quad (\text{G4})$$

$$y_{ic} = \begin{cases} 1 & \text{if class } c \text{ is the true class for sample } i, \\ 0 & \text{otherwise.} \end{cases} \quad (\text{G5})$$

We have implemented the quantum machine learning model using the PennyLane [94] library and use the JAX [95] backend to simulate it efficiently and train it. As an optimization paradigm, Adam is used in conjunction with a cosine decay schedule. Table VI lists the empirically chosen hyperparameters.

Hyperparameter	Value
learning rate α	1
n_{epoch}	200
n_{layer}	5
batch size	1000

Table VI. Hyperparameters for the QML experiment.

b. Evaluation metrics

We included more local attribution methods than introduced in the main text, we briefly present them here. As a starting point, one could simply take the local information present in the gradients. This leads to the explanation we refer to as `grad`: $E_i(x) = \partial_i f(x)$, and its averaged version with respect to local perturbations $E_i(x) = \mathbb{E}_{\xi \sim \mathcal{N}(0, \sigma^2 \mathbb{I})} [\partial_i f(x + \xi)]$, which we call `smooth_grad` [96]. `smooth_grad` was initially introduced to provide robustness of explanations against noise, as for instance

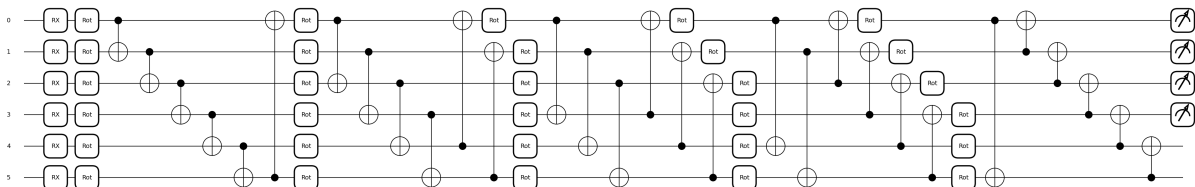


Figure 3. Schematic of the employed parametrized quantum circuit.

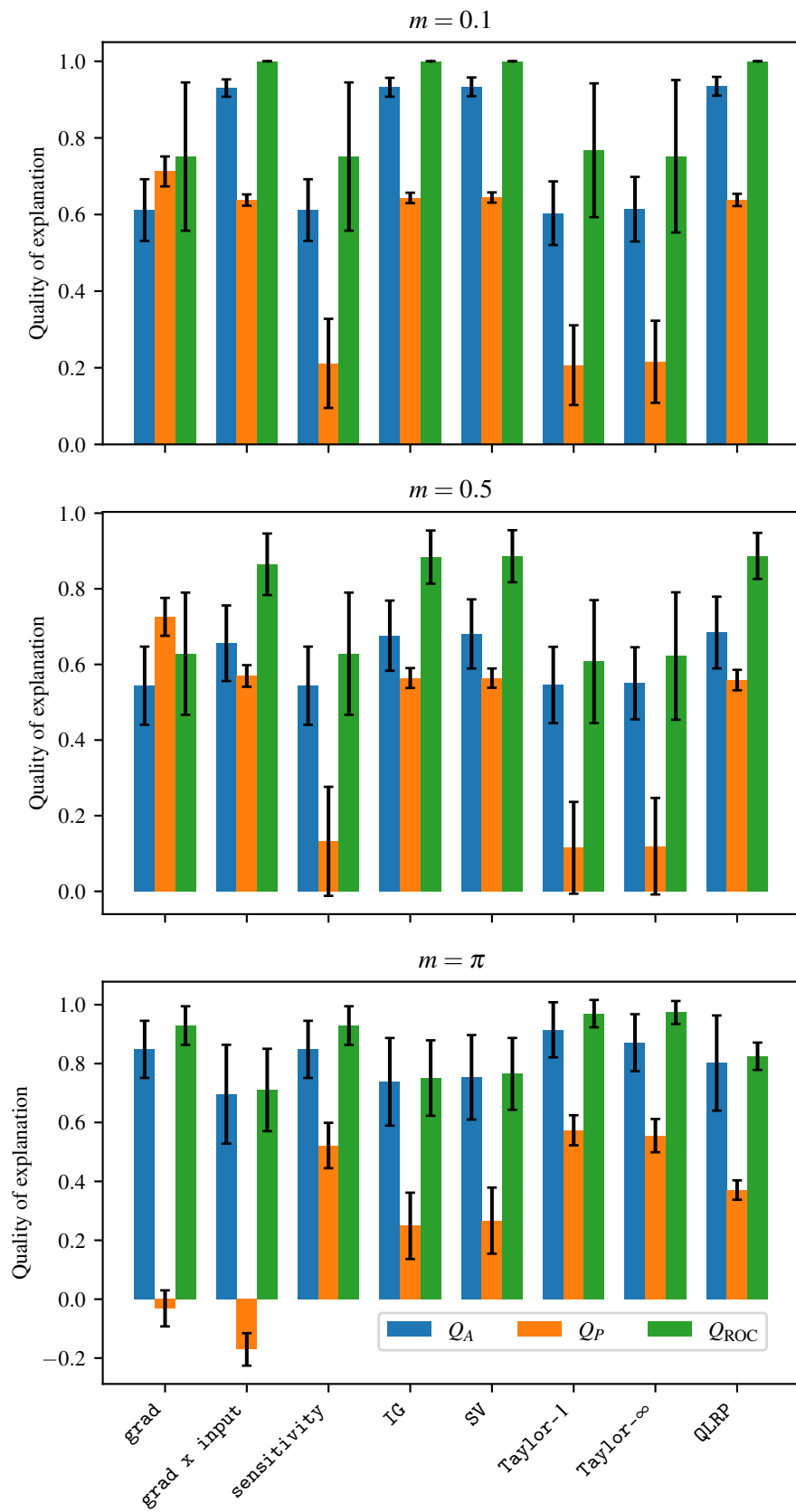


Figure 4. Full evaluation metrics

when the gradients of the model are discontinuous. In general neither `grad` nor `smooth_grad` give rise to positive explanations², which is addressed straightforwardly by the next explanation method, called `sensitivity`: $E_i(x) = |\partial_i f(x)|$ ³. Since the positivity of `sensitivity` comes from measuring the magnitude of the gradient, it may be the case that `sensitivity` does not relate to positive causation.

With the aim of comparing and evaluating the performance of the different attribution algorithms, we compute different metrics based on a desired ground-truth input mask and the local feature attributions across the data samples. To this end, we exploit the structure of the synthetic data-set. As discussed above in Appendix G 1, each of the four classes’ input data exhibits three important and three unimportant directions in \mathbb{R}^6 , by design. Thus, one can generate the ground truth masks for each class accordingly. The elements of the ground truth mask $M^c(x) \in \mathbb{R}^6$ for a sample $(\mathbf{x}_i, \mathbf{y}_i)$ belonging to class c is consequently given by

$$M_j^c(x_i) = \begin{cases} 1 & \text{if } j \in \{\text{main dimensions of class } c\}, \\ 0 & \text{otherwise.} \end{cases} \quad (\text{G6})$$

We compute the three evaluation metrics $Q_{\{A,P,ROC\}}$ for each input x_i , and then average over all inputs. The quality scores we report are thus $\mathbb{E}_x [Q_A(x)]$, $\mathbb{E}_x [Q_P(x)]$, and $Q_{ROC}(T)$, where T is a set of 200 inputs for each class.

Fig. 4 summarized our numerical results across all explanations and evaluation metrics.

² Requiring *positive* heat-maps $E_i(x) \geq 0$ might seem meaningful, as “positive causation” can be easier to interpret than “negative causation”. Yet, in practice that is too stringent an axiom to fulfill, so it is commonly not taken

into account.

³ Also standard is to call sensitivity the square of the gradient, instead of its magnitude.