

Evaluating Generalization in Classical and Quantum Generative Models

Kaitlin Gili,^{1,2,*} Marta Mauri,^{3,*} and Alejandro Perdomo-Ortiz^{3,†}

¹University of Oxford, Oxford OX1 2JD, United Kingdom

²Zapata Computing Canada Inc., 325 Front St W, Toronto, ON, M5V 2Y1, Canada

³Zapata Computing Canada Inc., 325 Front St W, Toronto, ON, M5V 2Y1

(Dated: May 27, 2022)

Defining and accurately measuring generalization in generative models remains an ongoing challenge and a topic of active research within the machine learning community. This is in contrast to discriminative models, where there is a clear definition of generalization, i.e., the model’s classification accuracy when faced with unseen data. In this work, we construct a simple and unambiguous approach to evaluate the generalization capabilities of generative models. Using the sample-based generalization metrics proposed here, any generative model, from state-of-the-art classical generative models such as GANs to quantum models such as Quantum Circuit Born Machines, can be evaluated on the same ground on a concrete well-defined framework. In contrast to other sample-based metrics for probing generalization, we leverage constrained optimization problems (e.g., cardinality constrained problems) and use these discrete datasets to define specific metrics capable of unambiguously measuring the quality of the samples and the model’s generalization capabilities for generating data beyond the training set but still within the valid solution space. Additionally, our metrics can diagnose trainability issues such as mode collapse and overfitting, as we illustrate when comparing GANs to quantum-inspired models built out of tensor networks. Our simulation results show that our quantum-inspired models have up to a $68\times$ enhancement in generating unseen unique and valid samples compared to GANs, and a ratio of 61:2 for generating samples with better quality than those observed in the training set. We foresee these metrics as valuable tools for rigorously defining practical quantum advantage in the domain of generative modeling.

INTRODUCTION

Utilizing generative models for unsupervised learning tasks has become increasingly popular within both the classical and the quantum machine learning (QML) communities (see e.g., [1–6]). With the development of deeper classical General Adversarial Networks (GANs) [7–9], quantum-inspired methods such as Tensor Network Born Machines (TNBMs) [10] and quantum models such as Quantum Circuit Born Machines (QCBMs) [4], generative models have been identified as leading candidates for advanced artificial intelligence tasks and quantum advantage applications. For example, generative neural networks can be used for molecular design and discovery [11, 12], image synthesis and deepfakes generation [13, 14], risk communication for malware defense [15], and secure modeling of private data across industries [16].

Recently, quantum and quantum-assisted generative models been introduced as a potential candidate for generative tasks ranging from the generation of images (see e.g., [17–19]) to the optimization of low-risk financial portfolios [20–22]. Despite the fact that these quantum models have been proven to learn distributions which are outside of classical reach [5, 23–26], we have yet to fully understand the power of these models and compare their generalization capabilities in the classical, quantum-assisted and hybrid quantum-classical regimes.

Indeed, in many applications, the real advantage behind machine learning (ML) is unlocking a trained model’s ability to generalize, where the generalization capability describes a model’s relationship with unseen data. The definition of

generalization takes on slightly different nuances depending on the ML task under examination, and as such, there have been many contributions that investigate this property given its remarkable importance [16, 27–30]. In the context of unsupervised generative learning, where models only have access to unlabelled data, generalization describes a trained model’s capability to generate quality data beyond the training set. This is different from the generalization error that exists in commonly used supervised discriminative models (e.g. image classification), which is equipped with a clear definition, even though its characterization is a popular and heavily investigated topic within the ML and QML community (see e.g., [28, 30–32]). Therefore, we restrict the scope of this work to investigating generalization in unsupervised generative models, which is far less discussed and understood, despite their growing impact and importance. In fact, even giving a quantitative and measurable definition of generalization, along with appropriate evaluation metrics, is far from being straightforward, as one can easily see from the wide variety of ideas that lack a unitary and well accepted vision [16, 33–37].

When evaluating unsupervised generative models, we tend to veer towards metrics that capture high quality and diversity among the generated data, but leave out novelty. Hence, with these metrics, generalization is never the focus of study. These are metrics that still provide a perfect score whenever the training dataset is exactly reproduced, a phenomenon known as data-copying [38]. In the case of classical and quantum models, metrics such as the Inception Score (IS) [39], Fréchet Inception Distance (FID) [40], and Precision and Recall [35–37, 41] are common examples of metrics that demonstrate the model’s ability to generate high quality and diverse data without providing a complete picture on how well the model can generalize. Other methods have attempted to touch on generalization by probing the model’s inductive bias, i.e., the

* Both authors contributed equally to this work.

† alejandro@zapatacomputing.com

intrinsic bias due to a predetermined architecture or training procedure that gives rise to generalization [27, 42]. However, this method fails to take into account the model’s ability to sift through poor quality unseen samples and generate only unseen data of value - in addition to the fact that a theoretical understanding of the exact role played by inductive bias within generalization is still an ongoing challenge in the classical and quantum machine learning community. Lastly, several works discussing quantum generative models have hinted at the concept of generalization, but have ultimately restricted their work to replicating a given target probability distribution, leaving such a question for future research [10, 20, 43–45]. A summary of previously proposed alternative metrics and evaluation schemes is presented in Appendix A.

If we stick to metrics that reward the model for reproducing the training set, then we are not allowing the model to learn what information, features, or patterns in the dataset are important and which are not - in essence, it doesn’t distinguish between what to keep and what to disregard. Without doing this, the model is simply memorizing. However, exactly as in human learning, the ability of ‘forgetting’ what is not important is a fundamental part of the learning process [46].

In this work, we present a unified approach to measure the generalization capabilities of both state-of-the-art classical and quantum generative models. Our proposal bridges the gap between two important research efforts. The first is the known rigorous and theoretical standpoint, which usually lacks immediate real-world application (see e.g., [26, 47]). Conversely, the other uses state-of-the-art real-world datasets, where only heuristics and approximate metrics can be proposed, but where the complexity of the models and tasks at hand blur any definite conclusions about the model generalization capacity. In our framework, leveraging discrete datasets relevant to many application domains [48], we can unequivocally measure the generalization capabilities of generative models. As in other sample-based metrics, our framework can be applied to any generative model, but the discrete nature of the sample space, combined with the unique “cost/value” associated with each sample, allows for an unambiguous quality assessment of the generated samples.

In Sections I and II, we first provide a robust definition of generalization, introducing concepts and a discrete dataset framework to describe and assess this capability. Building off of previous work [20], in Section III we introduce robust sample-based metrics that allow one to conduct a comprehensive quantitative assessment of a model’s generalization capabilities and detect common pitfalls associated to the training process. Furthermore, in Sections IV and V we illustrate our approach by comparing models from two separate regimes, namely fully classical GANs and quantum-inspired TNBM architectures, for a specific task with relevance in financial asset management.

To the best of our knowledge, this is the first proposal of an approach that combines a heuristic-based analysis with an application-based dataset to quantitatively evaluate generalization of unsupervised generative models, as well as the first direct comparison of generalization capabilities between classical and quantum-inspired ML models.

I. GENERALIZATION

Unsupervised generative models aim at capturing implicit correlations among unlabeled training data in order to generate samples with the same underlying features. In this work, we focus on binary encodings of datasets with discrete values, and therefore, discrete probability distributions. This is needed to facilitate the comparison of quantum and classical generative models, and to allow for a more accurate and unambiguous evaluation of generalization as opposed to the continuous case, as further clarified in Section IC.

More concretely, given a dataset $\mathcal{D}_{\text{Train}} = \{x_1, x_2, \dots, x_T\}$, where each sample x_t is an N -dimensional binary vector such that $x_t \in \{0, 1\}^N$ with $t = 1, 2, \dots, T$, we can train a generative model to resemble the unknown probability distribution $P(x)$ from which the samples in $\mathcal{D}_{\text{Train}}$ were drawn. Since a goal of the present work is to compare the generalization capabilities of models, we introduce formal definitions and metrics to quantify different aspects of the behaviours that arise when we sample from the generative model. We denote these samples as $\mathcal{D}_{\text{Gen}} = \{x_1, x_2, \dots, x_G\}$, where each x_g is again an N -dimensional binary bitstring, with $g = 1, 2, \dots, G$. As it will be shown later, the only requirement for the data distribution $P(x)$ is to have a support, which is a “valid” sector, and a complement, which is a set of noise or undesirable features. Many real-world datasets can be represented this way: for example, portfolio optimization as demonstrated in our work, as well as molecular design problems [49]. Remarkably, the notion of a constraint that defines valid and invalid spaces arises naturally within the context of combinatorial optimization as the constraint is usually part of the problem definition [48].

The formal definitions of our metrics are given in Section III, but prior to it, we provide a brief high level introduction to each of them, presenting the essential concepts for studying various flavours of generalization.

A. Pre-Generalization

We refer to *pre-generalization* as the generative model’s ability to go beyond the training set $\mathcal{D}_{\text{Train}}$ by producing unseen outputs. More precisely, for any level of generalization to occur it is necessary - but not sufficient - that there exist some points x_g such that

$$x_g \in \mathcal{D}_{\text{Gen}} \wedge x_g \notin \mathcal{D}_{\text{Train}}. \quad (1)$$

However, these outputs may not be samples distributed according to $P(x)$; for example, they may just be meaningless noise instead. In other words, pre-generalization is the model’s ability to generate any new output - whether it is distributed according to $P(x)$ or not (Figure 1). Note that we consider this behaviour to be a prerequisite for a model to be able to generalize, and not generalization in itself. As mentioned above and further specified below, to have any kind of generalization, a model must first be able to generate data beyond the training set, and the generalization potential is higher if the amount of unseen data is maximized. This implies that the training set cannot be exhaustive, i.e. the number

of unique¹ training bitstrings must be less than the number of unique bitstrings that can be sampled from $P(x)$. To discover new data, the training dataset should not consist of all of the bitstrings that could be sampled from the original distribution (i.e. its support).

The pre-generalization behaviour can be verified with our exploration metric E , defined in Section III A, that quantifies how many generated samples were not included in the training set. We note that this quantity has a similar definition to the *authenticity* metric in [16], that captures sample novelty. However, our exploration metric is computed directly from samples rather than requiring an embedding scheme and a separate classification network. This quantity allows one to investigate the general questions: “*Can the model reach out-of-training data points? And with which frequency?*”.

B. Validity-Based Generalization

We refer to *validity-based generalization* as the generative model’s ability to go beyond the training set $\mathcal{D}_{\text{Train}}$ and effectively produce new bitstrings living in a given solution space with the underlying distribution $P(x)$ (Figure 1). In other words, the model is able to learn a fixed particular feature about bitstrings drawn from $P(x)$ and produce new samples with the same feature, where this feature is specified via a constraint on the bitstrings. More precisely, the generative model outputs samples x_g such that

$$x_g \notin \mathcal{D}_{\text{Train}} \wedge x_g \in \text{support of } P(x). \quad (2)$$

We remark here that this approach for validity-based generalization is task-independent, as the metrics are exclusively sample-based and agnostic to the specific use case, or more specifically, independent of the quality associated to each bitstring. In Section II we highlight the essential conditions one needs to meet when defining an appropriate task to study validity-based generalization.

We evaluate the validity-based generalization behaviour introducing the three metrics of *fidelity* F , *rate* R , and *coverage* C . In a nutshell, F quantifies the probability that a model generates unseen samples that are valid results rather than unwanted noise. R quantifies the frequency at which a model produces unseen and valid results. C quantifies the fraction of unseen and valid results retrieved among all the potential valid and unseen samples. These metrics allow one to answer the following general questions, respectively linked to the three generalization estimators presented above:

- F : “How effectively can the model distinguish between noisy and valid unseen results?”
- R : “How efficiently can the model reach unseen and valid results?”
- C : “How effectively can the model reach all unseen and valid results?”

¹ Bitstrings = {00, 00, 11}, unique bitstrings = {00, 11}.

C. Quality-Based Generalization

We refer to *quality-based generalization* as the generative model’s ability to go beyond the training set $\mathcal{D}_{\text{Train}}$ and effectively produce bitstrings living in a given solution space with underlying distribution $P(x)$, where the new bitstrings can be mapped to a real number indicating their quality. While there can be many examples of functional maps that one could use to assign each bitstring a score to be maximized, we emphasize optimization as a natural choice for assigning such a value to each sample. In this case, the score is quantified by a cost to be minimized. In other words, optimization provides a natural framework to introduce quantitative estimators of generalization, as a generative task can be equipped with a well defined cost function, indicating the quality of samples. The framework presented here combines generalization and optimization as a promising strategy towards the definition of quantitative metrics.

When focusing on quality-based generalization, one is interested in generating samples that not only follow a given probability distribution of interest, but also have associated costs that minimize a given objective function (Figure 1). When considering continuous data distributions (e.g. in image generation tasks), assessing the quality of samples is particularly challenging, as embedding and non trivial transformations are needed in order to utilize the available metrics (see e.g., Refs. [16, 27]). Hence, on purpose we limit the scope of this work to discrete datasets, since this setting provides a more accurate and unambiguous evaluation of the generalization capabilities.

A generative model thus exhibits quality-based generalization if it is able to produce at least some unseen and valid samples that have on average similarly low (or lower) cost values than the ones associated to at least some of the training samples. More precisely,

$$x_g \text{ satisfies Eq. (2) } \wedge f(\mathcal{D}_{\text{Gen}}, C(x)) < f(\mathcal{D}_{\text{Train}}, C(x)), \quad (3)$$

for a given suitable function f (e.g., the minimum sample cost $C(x)$ in each sample set) that depends on how strict the cost minimization requirements are for the problem under examination (see Section III C).

Developing metrics for assessing quality-based generalization is a task-dependent challenge as it allows one to evaluate the model’s *sample quality*, according to a specific task and measured by its associated cost function.

In Section III C, we introduce two versions of the sample quality metric, induced by a different choice of f : the first one evaluates the model’s ability to generate a minimum cost value that is lower than anything in the training set, whereas the second accounts for a diversity of samples whose cost is below a user-defined percentile threshold. Even though the former could seem more adequate to quantify the generator’s ability to go beyond the sample quality available in the training set, it may be the case that producing the lowest cost value is not the only desired behaviour of the task. Furthermore, it may be that the desired behaviour is to generate diversity of new samples with a cost comparable to the lowest values

found in the training set. In this scenario, the latter version allows one to reward alternative solutions without restricting the model only toward values below the training threshold. Taking the sample quality metrics one-step further, we also see value in including the number of unique samples with a lower cost value than a user-defined threshold in the training set (e.g. the minimum value in the training set) in the evaluation process, as for many practical optimization tasks, one cares about reaching a diverse pool of these samples.

The quality-based generalization metrics allow us to investigate the general question: “*Can the model reach unseen and valid results that are more or just as valuable than the best in the training set?*”.

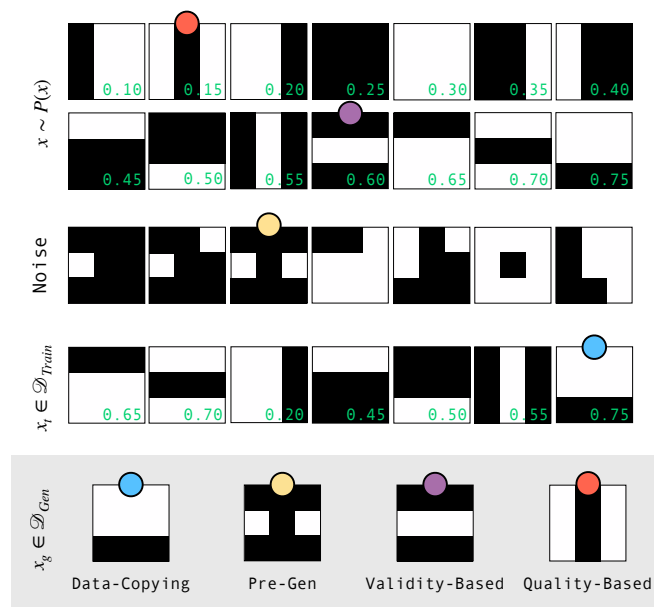


FIG. 1. **A visual representation of generalization-related concepts.** The figure shows the different behaviours a model can exhibit when generating data, using a 3x3 Bars and Stripes dataset as an example. The top two rows display a set of samples x distributed according to the data distribution $P(x)$; note that only a subset of the 3x3 Bars and Stripes dataset is displayed, rather than the full set of patterns. The Noise row contains samples that do not belong to this dataset. The fourth row contains a subset of samples $x_t \in \mathcal{D}_{\text{Train}}$ used for training and distributed according to $P_{\text{Train}}(x_t)$, while the bottom row shows a new set of samples x_g produced by the model and living in \mathcal{D}_{Gen} . Note that each sample contains an associated toy-score that corresponds to the samples’ associated cost. In this toy example, the samples are assigned a real-valued score in $(0, 1)$, except for noisy samples that don’t have an associated cost as they are not part of the valid solution space. The bottom row displays four samples from the generated queries, each of which is tagged with a different model behaviour: memorizing data from $\mathcal{D}_{\text{Train}}$ (blue dot), producing data outside of $\mathcal{D}_{\text{Train}}$ that may be noise (yellow dot), generalizing to new data distributed according to $P(x)$ (purple dot), and generalizing to new data distributed according to $P(x)$ that contains a minimum value to an associated cost function (red dot).

II. GENERALIZATION TASK DEFINITION

In order to properly assess generalization, the generative model’s task must meet some essential requirements. Such assumptions do not limit the scope of our approach as they simply provide a robust definition of the task at hand.

As previously specified, we focus our analysis on binary encodings of discrete datasets $\mathcal{D}_{\text{Train}} = \{x_1, x_2, \dots, x_T\}$, with $x_t \in \{0, 1\}^N$. We can thus identify a search space \mathcal{U} of size 2^N , that contains all possible N -dimensional bitstrings. For validity-based generalization, there must exist a subspace of \mathcal{U} containing the set of bitstrings we would like our trained model to generate. We refer to this as the valid solution space \mathcal{S} , that includes all the samples that exhibit a given desired feature. Hence, the model aims to approximate the underlying unknown *data distribution*, defined as:

$$P(x) = \frac{1}{|\mathcal{S}|}, \forall x \in \mathcal{S}. \quad (4)$$

We highlight that the notion of validity produces a non-trivial distribution of valid samples across the overall search space \mathcal{U} , adding complexity to the problem despite the data distribution being uniform over the solution space \mathcal{S} . We emphasize that this general solution space \mathcal{S} will contain different bitstrings for various representational datasets of interest. For instance, Figure 1 displays samples from the well known Bars and Stripes dataset [50]: in this case, the solution space \mathcal{S} would contain all valid bar and stripe patterns, some of which are shown in the top row of the figure. Alternative datasets could focus on solution spaces defined by a parity constraint, by a cardinality constraint or by any other property of interest. We highlight that the solution space must have a well defined notion of validity that can be evaluated for each of the bitstrings in \mathcal{U} to verify whether or not they live in its subset \mathcal{S} .

The model’s task is therefore to generate novel samples in \mathcal{S} , after a learning process involving a limited number T of unique training samples, i.e. $T = \epsilon|\mathcal{S}|$, where the seen portion $\epsilon \ll 1$ is a small parameter quantifying the percentage of \mathcal{S} that gets seen during training. Note that this is a necessary requirement for generalization because it guarantees that the training set is not exhaustive.

With T training samples, the model has access only to an approximated version of the data distribution, that we denote as the *training distribution*:

$$P_{\text{Train}}(x) = \frac{1}{T}, \forall x \in \mathcal{D}_{\text{Train}}. \quad (5)$$

For quality-based generalization, there is an additional requirement as this behaviour depends not only on the validity of the bitstrings, but also on the value associated to each pattern, according to a cost function $C(x)$. As such, in order to assess quality-based generalization, it is necessary for the task of interest to have a well-defined objective function that indicates the cost of each bitstring, in search for minimum values.

As we would like for our model to learn the valid bitstring patterns as well as to generate patterns with low-cost values, it

is integral to bias the dataset distribution in Eq. (4). Here we use a *softmax* function in order to introduce cost-related information in the training data set. In this scenario, the training samples approximate the following *biased training distribution*:

$$P_{\text{Train}}^{(b)}(x) = \frac{e^{-\beta_m C(x)}}{\sum_{i=1}^T e^{-\beta_m C(x)}}, \forall x \in \mathcal{D}_{\text{Train}}. \quad (6)$$

Following Ref. [20], $\frac{1}{\beta_m}$ was chosen to be the standard deviation of the costs in the training data, whereas $C(x)$ is the cost of each sample bitstring [20].

In summary, the two main essentials for evaluating validity-based and quality-based generalization, respectively are the following:

- There exists a well-defined solution space \mathcal{S} , containing bitstring patterns that are valid according to easy to specify and verify constraints.
- There exists a well-defined cost function $C(x)$ that can be computed to assess the generalization for all valid bitstring patterns.

III. METRICS FOR EVALUATING GENERALIZATION

As described in Section I and Section II, generalization occurs when we retrieve novel samples that have features corresponding to some underlying distribution, approximated by the trained generative model. To achieve this, typical training schemes attempt to minimize the distance between the training and the generated distribution using a similarity measurement such as the Kullback-Leibler (KL) divergence. Oftentimes, these divergences are further utilized to evaluate the quality of the model after training [51]. However, a necessary and sufficient condition for the KL divergence to be exactly zero is that the training and the generated distributions are identical, namely:

$$\text{KL} = 0 \Leftrightarrow P_{\text{Train}} = P_{\text{model}},$$

which coincides with pure memorization. Hence, by using this metric for evaluation purposes, we are nudging the generalization analysis to focus only on the data-copying capability. Therefore, a KL divergence-like measure is not suitable in the context of generalization as its optimal value does not necessarily imply good generalization behaviours, since $P_{\text{Train}}(x)$ is different from the desired $P(x)$.

Alternatively, another pair of metrics that has been proposed for evaluating training and generalization performances of generative models [34, 35] are the so-called precision and recall. Precision p quantifies the model's ability to produce samples in the solution set \mathcal{S} , and recall r measures the model's ability to recover all solutions in \mathcal{S} [16]. However, these metrics are not an ideal choice in the context of generalization as they do not necessarily focus on unseen samples. In other words, the definitions of precision and recall do not distinguish between queries that are in the training set from

those that are not, so much so that when one is interested in assessing whether a query is a novel data point, they need to introduce other evaluation tools, such as authenticity [16] to take this feature into account. As an extreme example, if the training set is exhaustive and all samples are retrieved from the discrete dataset, the recall of the model would be perfect although no generalization occurred. We introduce a deeper discussion of various metrics used for generative model evaluation in Section A.

To improve upon these metrics that do not fully capture generalization, we propose a novel approach for evaluating the generalization capabilities of quantum and classical generative models. To give a quantitative definition of the generalization metrics, we first need to clarify the nomenclature of all the spaces involved. We have already defined the set of all queries retrieved from a trained generative model as \mathcal{D}_{Gen} , where $|\mathcal{D}_{\text{Gen}}| = Q$. We then call \mathcal{G}_{sol} the set of all valid and unseen non-unique queries, which reflect the model's validity-based generalization capability. We further define a subset of \mathcal{G}_{sol} that contains all its unique bitstring solutions as g_{sol} , thus the only difference between \mathcal{G}_{sol} and g_{sol} is that in the latter each bitstring appears only once, whereas in the former there can be many occurrences of the same sample. Lastly, we define the subset of unseen non-unique queries as \mathcal{G}_{new} , where some of these queries might be unwanted noise and hence reflect the model's exploration capability. Note that we use uppercase variables for non-unique sets and lower case variables for unique sets, and a visual representation of the sets in play can be found in Figure 2.

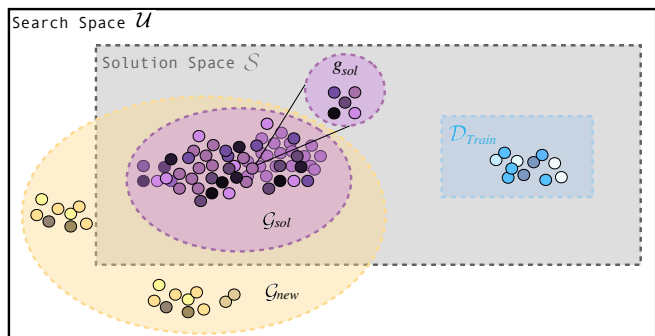


FIG. 2. **A visual representation of all possible spaces where a generated query might be located.** Each query is represented by a color-coded dot, where the color-code is the same as in Figure 1 (Data-Copying: blue, Pre-Generalization: yellow, Validity-Based Generalization: purple) and the color-shade represents a unique bitstring sample. We take all non-unique queries outside of the training set to be in the subset \mathcal{G}_{new} (inside the yellow oval), whether they are in the solution space \mathcal{S} or not. Furthermore, we take \mathcal{G}_{sol} to be all non-unique queries that exist in the solution space (inside the pink oval) and g_{sol} to be all of the unique queries among \mathcal{G}_{sol} (zoomed-in). Lastly, if a query exists in $\mathcal{D}_{\text{Train}}$, it is a memorized count from the training set. We note that the quality-based queries (not shown) must exist inside of the solution space.

Having clarified the nomenclature of the spaces involved in the task, we can now proceed to the definition of the generalization metrics.

A. Evaluating Pre-Generalization

While a model’s capability to generate unseen samples that are not valid or valuable solutions to the task at hand is not considered generalization behaviour in itself, it is an important prerequisite for generalization. If the model is not able to go beyond the training set, even just to produce noisy outputs, then the model is not passing the first requirement for generalization - the ability to produce novel data points. To conduct a pre-generalization evaluation prior to assessing for any kind of validity-based or quality-based generalization, we introduce the *exploration* metric E , that quantifies the fraction of generated queries that are new data points, namely:

$$E = \frac{|\mathcal{G}_{\text{new}}|}{Q}. \quad (7)$$

If $E \approx 0$, the model will not pass the first required check for generalization. This may be due to an intrinsic property of the model, i.e., the inability to generate novel data, or it can be an artifact of the training set being (almost) exhaustive, because nothing new can be generated if the training data covers (almost) all the entire valid space.

B. Evaluating Validity-Based Generalization

We introduce three sample-based metrics that describe each model’s validity-based generalization behaviour after training: fidelity F , rate R , and coverage C .

Fidelity describes the model’s ability to distinguish an unseen and valid sample in \mathcal{S} from a meaningless output (i.e., noise) and it quantifies the fraction of unseen queries that fall into the unseen solution space. It is defined as follows:

$$F = \frac{|\mathcal{G}_{\text{sol}}|}{|\mathcal{G}_{\text{new}}|}. \quad (8)$$

Rate describes the model’s ability to efficiently produce unseen and valid samples and it quantifies the fraction of all queries that fall into the unseen solution space, namely:

$$R = \frac{|\mathcal{G}_{\text{sol}}|}{Q}. \quad (9)$$

Coverage describes the model’s ability to recover all unique unseen and valid samples and it quantifies how much of the solution space that was unexplored gets covered by the generative model’s queries. It is defined as follows, where we highlight that the ratio does not take into account the queries’ frequencies, as a single occurrence has the same weight as a one that appears multiple times:

$$C = \frac{|g_{\text{sol}}|}{|\mathcal{S}| - T}. \quad (10)$$

We highlight that one should expect the value of these metrics to depend on the number of queries Q that are retrieved from the trained model. For example, to have a good coverage of a space, i.e., $C \rightarrow 1$, one should have enough samples

that fall in the entire unexplored space. However, this dependency does not constitute a limitation for drawing a comparison between models, as we can fix the number of queries for all the models under investigation, and evaluate and fairly compare their generalization performance at the given number of queries. Moreover, in Section V B, we further showcase the values of C as we increase the number of queries toward and beyond the size of the solution space. We see a clear trend towards the metric ideal limit $C \rightarrow 1$ as we increase the number of queries. Conversely, in Appendix E we demonstrate that fidelity and rate are not dependent on the number of generated samples, despite being sample-based metrics.

We note that the different metrics are not completely independent, as there are mutual relations between them. For instance, it can be noted that rate and fidelity are correlated, as $R = EF$. Rate is the same as fidelity whenever a model generates exclusively unseen queries, which only holds in the case of perfect generalization (or in pathological cases such as mode collapse to unseen and valid queries). Another example of mutual relation between the metrics is that $C \leq \frac{EQ}{|\mathcal{S}-T|}$, which implies that $C < E$ for large solution spaces and limited queries budget.

To further clarify the expected metrics for a well-generalizing model, we highlight that these metrics will be exactly 1 when evaluated for a model that exhibits the highest validity-based generalization. However, in a practical sense, this might be difficult to achieve; we are then equipped with a theoretical unitary upper bound for all metrics, with the understanding that one should aim to reach toward this limit to obtain a robust model for generalization.

Lastly, we note that the pre-generalization condition in Eq. (1) impacts the validity metrics, and hence, exploration E is directly related to (F, R, C) . For F , the pre-generalization condition in Eq. (1) must be met in order for the metric to be well-defined. When the condition is not met, F will be null, and $C, R = 0$. Therefore, our metrics rely on the model’s ability to go beyond the training set, and will indicate if the model is only data-copying. Other properties from the model can be inferred from these metrics as demonstrated in Table IV in Appendix C. For example, a metric which measures the degree of data-copying could be defined as $D = 1 - E$, hence perfect memorization would mean $E = 0$. We highlight that, in this framework, one can additionally use our proposed metrics to detect alternative and complementary behaviours to generalization and define additional metrics that are tailored towards specific properties one would like to investigate.

In conclusion, we propose to utilize the metrics (F, R, C) to introduce a 3D quantitative investigation of the generalization capabilities mentioned in Section I B, that we report here for convenience:

- Fidelity, F , evaluates how effectively the model can distinguish between unseen valid and invalid bitstrings.
- Rate, R , evaluates how efficiently the model can produce unseen and valid bitstrings.
- Coverage, C , evaluates how effectively the model can retrieve all unseen and valid patterns.

C. Evaluating Quality-Based Generalization

To quantify the quality-based generalization properties of a generative model, we propose adequate metrics addressing the *sample quality* of the generated samples, which speak to how many of the queries are more valuable results in the context of a specific application domain, i.e., how many bitstrings have a low enough associated cost. Since the quality of a result depends on a given cost function, this metric is task-specific, as opposed to the validity-based generalization case that only requires the notion of validity of a query, according to a well defined hard constraint.

More precisely, we introduce different nuances of this *sample quality* metric for our quality-based generalization assessment, proposing two different versions with slightly different implementations of f in the right-hand side condition of Eq. (3).

Firstly, we consider the Minimum Value (*MV*) of the costs associated to the queries generated by the model as a relevant evaluation metric, since in many optimization applications the main goal is to find the solution that minimizes the cost, or equivalently, the sample with the best quality. This corresponds to choosing $f = \min$, so that the condition of Eq. (3) becomes:

$$x_g \text{ satisfies Eq. (2)} \wedge \min_{x_g \in \mathcal{D}_{\text{Gen}}} C(x_g) < \min_{x_t \in \mathcal{D}_{\text{Train}}} C(x_t). \quad (11)$$

Despite its practical impact, this punctual metric can be highly unstable if it is not supported by enough statistics as the metric relies on generating one specific value, the lowest. Since generating the query with the lowest cost is highly dependent on the selected batch b of queries, we define this metric as an average across B batches of queries to avoid biasing the results due to an anomalous batch. In other words,

$$MV = \frac{1}{B} \sum_{b=1}^B \min_{x_g \in \mathcal{G}_{\text{sol}}^b} C(x_g)$$

for each generative model being evaluated. For the results presented in this work, we fixed $B = 5$. Including such average in the definition of the *MV* metric itself contributes to alleviate its intrinsic instability, thus making it a more robust metric for quality-based generalization evaluation.

Secondly, we define the Utility U as the average cost of a user-defined set P_t of unseen and valid samples from the generative model. Specifically, $P_t(\mathcal{D})$ is the set obtained from taking the $t\%$ of samples with the best quality (lowest costs) in \mathcal{D} . Setting $t = 5$, this corresponds to choosing $f = \langle \cdot \rangle$ on the set P_5 , and the condition of Eq. (3) reads:

$$x_g \text{ satisfies Eq. (2)} \wedge \langle C(x_g) \rangle_{x_g \in P_5(\mathcal{G}_{\text{sol}})} < \langle C(x_t) \rangle_{x_t \in P_5(\mathcal{D}_{\text{Train}})}. \quad (12)$$

Given its set-based definition, this metric is much more stable than the previous one.

Lastly, we note that it is possible to give another definition of *sample quality*, which simply consists in counting the

number of unseen and valid queries whose cost is lower than a specific critical cost value $C'(x)$ in the training set. For example, one could take $C'(x)$ to be the lowest cost value in the training set i.e., $C'(x_t) = \min_{x_t \in \mathcal{D}_{\text{Train}}} C(x_t)$. When utilizing this estimator, one is interested in verifying the following condition:

$$\left| \{x_g \text{ s.t. } C(x_g) < C'(x_t)\} \right| > 0, \text{ for } x_t \in \mathcal{D}_{\text{Train}}, \quad (13)$$

where clearly a higher value of the left-hand side implies a better sample quality. Even though this quantity can carry interesting information, we don't include it among our quality-based generalization metrics as it is a harsh restriction to impose and may only be important for optimization tasks that are looking for many potential *MV* bitstrings. We highlight that our framework is not limited to the metrics proposed so far, but allows one to define several other figures of merit which can be relevant for specific applications at hand.

We use these metrics to introduce insight into a model's quality-based generalization capabilities, and determine which models are able to generate the most value for task-specific challenges. We emphasize again that this approach can be utilized beyond cost minimization problems, as long as there is a quantitative quality scale associated to each bitstring in the valid subspace.

IV. APPROACH DEMONSTRATION

To present the robustness of our approach in evaluating and comparing generative models, we choose a well-defined task and two families of models: classical GANs and quantum-inspired TNBMs. The following sections outline the specific use case (Section IV A) and the generative models (Section IV B) selected for our experimental demonstrations.

A. Use Case

To demonstrate a practical application of our approach, we choose an important use case in the finance sector that addresses the challenge of cardinality-constrained portfolio optimization. The goal of such task is to minimize the risk σ associated to a collection of assets, randomly selected from the S&P500 market index, for a fixed desired return ρ . Below, we highlight how this task is amenable to the framework and requirements described in Section II.

Given a fixed size N of the asset universe, a portfolio candidate can be encoded into a bitstring of length N , where each bit corresponds to an asset either being selected in the portfolio (1) or left out of the portfolio (0). Therefore, the search space \mathcal{U} of all possible portfolios grows exponentially with the asset universe size, i.e. $|\mathcal{U}| = 2^N$.

To assess validity-based generalization within this task, we define the solution space \mathcal{S} to be comprised of all bitstrings containing a fixed number $k = \frac{N}{2}$ of selected assets, i.e. a candidate solution must be a bitstring with a fixed Hamming weight equal to k .

With such k -cardinality constraint, the problem solution set \mathcal{S} contains all possible portfolio bitstrings x that fit this constraint. Thus, its cardinality is:

$$|\mathcal{S}| = \binom{N}{k}. \quad (14)$$

To further assess quality-based generalization, we define an objective function that encodes the quality of each bitstring, namely the financial risk σ associated to each portfolio, which in the case of the Mean-Variance Markowitz model [52] can be efficiently computed by means of Mixed Integer Quadratic Programming (MIQP) [21]. Unlike when investigating validity-based generalization, we use σ to bias the training dataset with the softmax function described in Eq. (6).

As such, this task satisfies both conditions necessary in order to evaluate validity-based and quality-based generalization. We again emphasize that our framework can be applied to any task that meets the essential requirements in Section II, and is not limited to this financial application.

B. Generative Models

We focus our investigation on Generative Adversarial Networks (GANs) and Tensor Network Born Machines (TNBMs). This choice is motivated by several reasons. On the one hand, GANs constitute one of the most popular and top utilized classical generative models, notwithstanding the challenges that plague their training such as mode collapse [53], convergence issues [54], and vanishing gradients [55]. Moreover, they are made up of several components that can be independently and successfully promoted to a quantum model [19], thus paving the way to the study of hybrid quantum-classical generative models. On the other hand, recent results for training TNBM architectures show that such model is a promising candidate to exhibit both validity-based and quality-based generalization behaviours [20]. We started our generalization study choosing these two models, but our approach can be leveraged to characterize any other state-of-the-art generative model of interest, and we do hope other interesting works will spin out from this initial proposal to evaluate quantitatively their generalization power. Future work can include an analysis of fully quantum models, even trained on hardware, once current limitations in training large and deep circuits are overcome.

1. Generative Adversarial Network (GAN)

Our classical model consists of a Generative Adversarial Network (GAN) architecture with a normal prior distribution, and we conduct the training as typically described in the literature [7–9]. GANs are trained as two neural networks, a discriminator D and a generator G , competing against one another for optimal performance in an adversarial game. Samples from a prior distribution $q(z)$ are fed into the generator's

input layer, and throughout training the generator attempts to produce new data x that can fool the discriminator into classifying x as a real rather than an artificially created data point. The goal of training is to maximize the generator's score and minimize the discriminator's score as described by the loss function:

$$C_{\text{GAN}} = \min_G \max_D [\mathbf{E}_{x \sim P_{\text{Train}}}(x) [\log D(x)] + \mathbf{E}_{z \sim q(z)} [\log(1 - D(G(z)))]]. \quad (15)$$

For both the generator and the discriminator, we utilize a feed forward architecture with fully connected linear layers (details are listed in Table III in Appendix B).

2. Tensor Network Born Machine (TNBM)

Our quantum-inspired generative model is a Tensor Network Born Machine (TNBM), whose underlying architecture is chosen to be a Matrix Product State (MPS), a well-known 1D tensor network characterized by a low level of entanglement [10]. A TNBM takes unlabelled N -dimensional training bitstrings from the dataset $\{x_t\}_{t=1}^T$, and aims to encode the underlying probability distribution in a quantum wavefunction ψ , expressing the correlations between samples in the amplitude of a quantum state, namely:

$$|\psi\rangle = \sum_{\{s\}} \sum_{\{\alpha\}} A_{\alpha_1}^{s_1} A_{\alpha_1 \alpha_2}^{s_2} \dots A_{\alpha_N}^{s_N} |s_1 s_2 \dots s_N\rangle. \quad (16)$$

To motivate this representation, we note that an N -dimensional bitstring can be interpreted as a possible realization of the spin state (0,1) of N particles $|s_1 s_2 \dots s_N\rangle$, and therefore the full quantum state can be written as a superposition of all the possible spin states. Rather than using the exact coefficient matrix to build $|\psi\rangle$, we approximate it by the product of smaller parametrized single-particle matrices A^{s_i} , where the dimensions $\{\alpha\}$ are known as bond dimensions. The summation across α determines the probability amplitude for each superposition state of individual sites; thus, the bond dimensions controls the expressivity of the TNBM.

We use a similar training method as described in [10], where models are trained via a DMRG-like algorithm with the log-likelihood cost function:

$$C(\theta) = -\frac{1}{T} \sum_t \log(p_\theta(x_t)). \quad (17)$$

During training, samples are generated from the wavefunction according to the Born Rule:

$$p_\theta(x_t) = |\langle x_t | \psi \rangle|^2, \quad (18)$$

and the goal of the learning process is to find an optimal TNBM parametrization θ such that $p_\theta(x_t) \rightarrow P_{\text{Train}}(x_t)$.

A TNBM is known as a quantum-inspired technique as it builds upon fundamental concepts and formalism of the quantum-mechanical theory, but it is executed entirely on a classical platform.

V. RESULTS AND DISCUSSION

Having defined several quantitative metrics that allow one to conduct a generalization analysis of generative models, we use them to investigate the performance of TNBM and GAN architectures. We present the results of our simulations, whose details are specified in Section V A. We demonstrate the robustness of our proposed metrics (Section V B), show their ability to spot common pitfalls in model training (Section V C), and introduce insights into the validity-based and quality-based generalization capabilities of each model (Section V D).

A. Simulation Details

For our experiments, we consider a specific instance of a cardinality constrained portfolio optimization task, where we aim at minimizing the associated risk σ for a given target return $\rho = 0.002$, such that the asset universe from which one can pick to build a new candidate portfolio has size $N = 20$. Here, assets are randomly selected from the S&P500 index, as previously done in [20, 21], and the return level ρ is the same as used in previous studies. We impose the cardinality constraint that each portfolio must have a fixed Hamming weight $k = \frac{N}{2} = 10$. As previously stated, such an essential restriction creates a subset of the search space \mathcal{U} , of size $2^N \sim O(10^6)$, defining a solution space \mathcal{S} of size $\binom{N}{k} \sim O(10^5)$. The choice of these values allows for a big enough space so that generalization capabilities can be probed.

Given the solution space of portfolio candidates, the data distribution $P(x)$ given in Eq. (4) used to assess validity-based generalization is automatically defined. To build a non-exhaustive $P_{\text{Train}}(x)$ as in Eq. (5), only a fixed number $T = \epsilon|\mathcal{S}|$ of training samples are randomly selected from the solution space, thus making the task of learning the distribution $P(x)$ a highly non-trivial (despite it being defined as a uniform distribution over the valid bitstrings). Specifically, all generative models are trained for a fixed number of epochs $n_{\text{epochs}} = 100$ with a fixed value of T that equals 1% of the solution space (i.e. $\epsilon = 0.01$), leaving the remaining 99% of the space available for testing their generalization capabilities. Several values of this hyperparameter have been investigated, and we found this particular percentage to be a good choice as it gives the models many chances of generalizing, while providing enough samples $T \sim O(10^3)$ for the learning process to be successful. In order to assess quality-based generalization, we conduct the same process outlined above, with the addition of a pre-processing step that uses a softmax function to introduce risk-based information in the training dataset, so that low-risk portfolios are assigned a higher probability, and sampled with higher frequency.

We investigated the generalization behaviours of different versions of the TNBM and GAN architectures, using various hyperparameter sets. In the case of the TNBM, we considered different values for the bond dimension α , as this is the main parameter that affects the model quality. For GANs, the

choice of hyperparameters is significantly more challenging [56]. Therefore, in addition to identifying hyperparameters via a trial-and-error procedure, we investigated whether automated hyperparameter optimization using Optuna [57] could significantly improve the performance. We propose three different GANs that only differ in their hyperparameters as per Table III in Appendix B, and show generalization behaviours for all of them. From here onward, we refer to a GAN that has mode collapsed onto one seen and valid bitstring as GAN-MC and to the Optuna enhanced GAN as GAN+.

As mentioned above, all models have been trained for a fixed number of epochs and the associated generalization metrics have been computed based on a fixed number $Q = 10^5$ of queries retrieved from the trained model returned after the last epoch. Other strategies can be employed, such as considering the set of weights associated to the lowest loss function during training, or including more advanced training techniques such as early stopping. We decided to leverage a simple training scheme to avoid introducing any training bias and allow for the fairest comparison of the two models under examination. We also chose to sample this high magnitude of queries since this was not a limitation for the problem size considered here. However, in Appendix E we present the behaviour of our sample-based metrics as a function of the number of queries.

All of the numerical experiments in this work were carried out with Orquestra[®] ² for workflow and data management.

B. Metric Robustness

The first step to validate our approach consists in showing the robustness of our sample-based metrics. To verify this, we conduct a statistical analysis of the generalization metrics' values and investigate the statistical errors associated to them.

We focus the robustness analysis on one instance of each of the two generative models presented in Section IV. Specifically, we consider a TNBM model with fixed bond dimension $\alpha = 7$, which has proven to be a good choice for generalization purposes as will be explained in Section V C. For GAN, we consider the set of hyperparameters displayed in the first column of Table III in Appendix B, which were selected as reasonable values via a trial and error procedure (i.e. without leveraging automated hyperparameter optimization). The analysis can be extended to other instances to further strengthen the evidence of the robustness of our metrics.

After training these two model instances using gradient-based optimizers (see Table III), we perform 30 independent query retrievals and compute our generalization metrics on these distinct sample sets. We then evaluate the relative percentage error³ associated to each of the metrics to assess their statistical robustness. For each of the two models, the error values for both validity-based and quality-based metrics are shown in Figure 3.

² <https://www.orquestra.io/>

³ Relative percentage error is defined as the standard deviation of the metric values over their average.

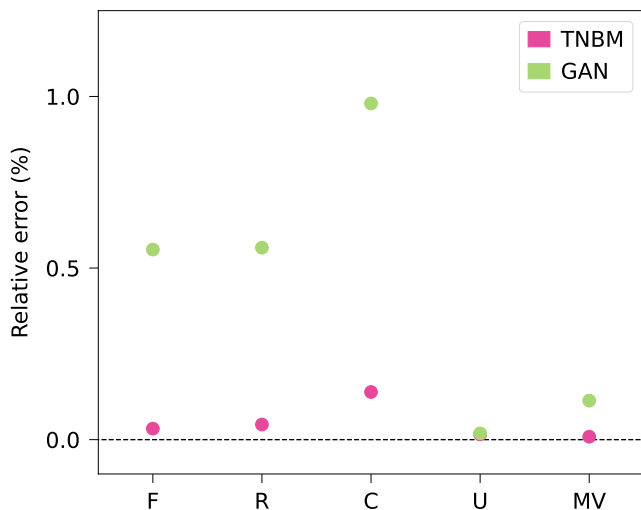


FIG. 3. **Robustness of the generalization metrics.** The plot shows the relative percentage error associated to each of the generalization metrics proposed in Section III, listed on the x axis. The errors are estimated as the relative standard deviation of independent metric values computed on 30 sets of queries generated by trained TNBM (pink) and GAN (green) models. The proposed metrics show their statistical robustness: the associated error is small, suggesting that our approach is sample-based but not sample-dependent. Henceforth, new independent same-size query batches from the trained model will produce similar metric results.

The errors associated to the different metrics assume similar values for the TNBM and GAN: this supports our claim that our metrics are model-agnostic and can be used to evaluate generalization capabilities for any generative model of interest. Furthermore, we can see in Figure 3 that the relative errors are less than 1%, thus suggesting that our metrics show significant robustness when computed on different sets of queries. Hence, we can affirm that the metrics proposed in this work are sample-based but not sample-dependent across different query batches of the same size.

The latter statement requires further clarification in the case of the coverage metric in Eq. (10). In this case, even though the coverage does not depend on the set of queries, it does depend on the number of queries that are retrieved from the trained model, as suggested in Section III B. The ideal unitary value of coverage is reached in the limit of a large number of queries, when the trained model has the opportunity to generate enough samples to cover most of the solution space. However, we note that, given a query budget Q , the effective upper bound UB to the coverage value is set by

$$UB = \frac{\min(Q, |\mathcal{S}|)}{|\mathcal{S}|} \leq 1,$$

thus implying that the ideal unitary value can be reached only with a sufficiently high number of queries, i.e. $Q \geq |\mathcal{S}|$. We investigated if the models considered so far show this trend as we increase the number of queries retrieved after training from 10^4 to $3 \cdot 10^6$. The results of the simulations are displayed in Figure 4; in Appendix D we compare them with the baseline

given by random sampling from the search space \mathcal{U} . Results for how the other metrics vary with the number of queries Q are shown in Appendix E.

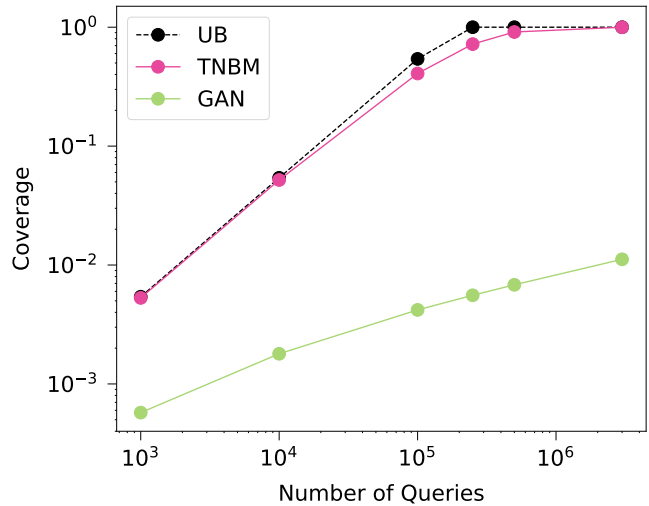


FIG. 4. **Coverage trends for increasing number of queries.** The plot displays the behaviours of the coverage metric for both TNBM (pink) and GAN (green) as we increase the number of queries Q retrieved from the trained models. The dashed black line shows the upper bound UB for each number of queries selected - i.e., the number of queries selected over the total size of the solution space. In the case of the TNBM, we observe that the coverage value follows the UB curve and saturates to the ideal value of 1.0 for large numbers of Q , corresponding to the scenario in which the trained model is able to generate all unseen and valid samples. In the case of the GAN, we still observe that the coverage value gets closer to UB and the ideal unitary threshold when more and more queries are drawn from the model. However, it remains further from UB and never reaches the desired threshold, suggesting that our GAN requires more queries than the TNBM to be able to reach all the unseen samples in the solution space.

The data shows that the TNBM coverage closely resembles the UB trend for any given value of Q and saturates to the ideal value of 1 for a large enough number of queries, implying that this model is able to achieve excellent coverage. Conversely, the GAN coverage is further from the UB and slowly increases without getting to the desired threshold, thus suggesting that significantly more queries would need to be taken to achieve a perfect exploration of all the unseen and valid patterns. Since there is no guarantee that the desired threshold is reached with a finite number of queries, this result might as well indicate that the model is quite poor at generalizing due to a high number of unreachable patterns. This is particularly relevant in the case of very large solution spaces \mathcal{S} . In this circumstance, the coverage metric has an intrinsic limitation: its low value might indicate that the number of generated queries is insufficient ($Q \ll |\mathcal{S}| - T$), rather than being due to poor generalization ($|g_{\text{sol}}| \approx 0$). Therefore, in order to mitigate the above issue when evaluating single models in the case of large problem sizes, we envision the denominator in C to be replaced by the number of queries Q . This solu-

tion will slightly distort the meaning of coverage in Eq. (10) to a new metric quantifying the rate at which the model generates unique unseen and valid samples. When extending to large problem sizes, we see this as a more relevant evaluation metric as one cares more about the diversity of unique unseen and valid samples the model can reach rather than reaching all of them, which would be impossible without the number of queries being at least the size of the solution space. However, as our experiments are conducted with a mid-sized problem space, we stick to the definition in Eq. (10) for our evaluation.

Even though the coverage metric is dependent on the number of queries and its interpretation in terms of generalization is affected by the size of the solution space, we can draw a fair comparison between the coverage of different models. Indeed, we can compare TNBM and GAN models if we keep the number of queries generated from each fixed, as reported in Section VD, where it will be shown that the quantum-inspired model outperforms this GAN model when given the same sample budget.

C. Spotting Pitfalls in Generative Model Training

We further demonstrate that we can use our metrics to detect common pitfalls that are known to affect the training of the TNBM and GAN models. This result strengthens the validity of our approach, which turns out not only to be a good framework for quantifying generalization of generative models, but also to enhance the study of their trainability. In the following sections, we show an example of this study for each of the models. For the TNBM, we analyze the relation between the bond dimension α , our generalization metrics, and the trainability of the model. Conversely, for the GAN, we investigate the relation between our metrics and mode collapse. Additional results to compare the training stability of the two classes of models are shown in Figure 17 in Appendix B.

1. TNBM Bond Dimension and Trainability

In the TNBM architecture, the bond dimension α of the MPS plays an important role in the model’s ability to generate good quality samples as it is directly correlated with the expressivity of the model. Typically, increasing the bond dimension leads to a better model approximation. We take this one step further and directly connect bond dimension to the model’s generalization behaviour and trainability.

In light of this goal, we train five different instances of the TNBM architecture on a fixed training dataset with various bond dimensions $\alpha \in \{3, 5, 7, 9, 11\}$. For a given α value, we select a typical⁴ training and build a model with the last set of parameters retrieved after the learning process. We

⁴ A typical training instance is identified as the resulting model from the median value of the loss function (i.e. KL divergence) at the last epoch, out of 30 independent trainings.

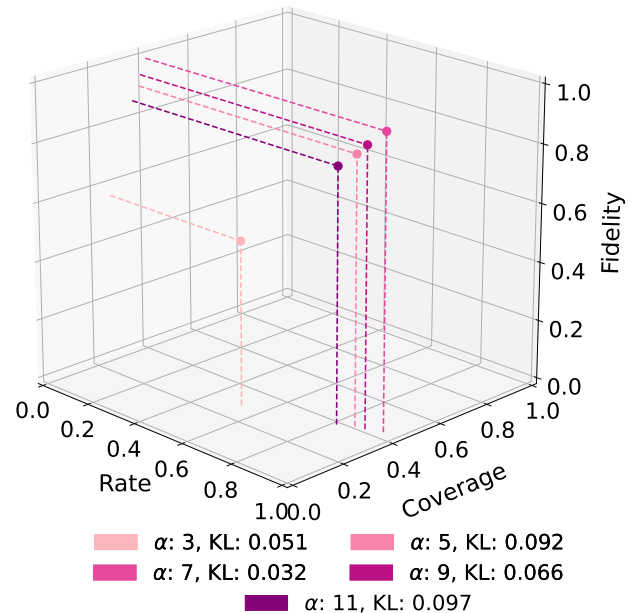


FIG. 5. **Training and generalization behaviours of the TNBM with different bond dimensions.** The plot displays the 3D evaluation of the validity-based generalization capabilities of the TNBM models with various $\alpha \in \{3, 5, 7, 9, 11\}$. Each data point corresponds to the average metrics’ values, whose associated error is too small to be visible on the plot. The legend connects each α to the last KL divergence value in the training after 100 epochs. The plot demonstrates that choosing $\alpha = 7$ provides the lowest KL divergence value along with the best generalization performance, thus establishing a link between this capability and trainability properties of generative models.

then generate 15 independent query batches from the trained model and compute our validity-based generalization metrics (F, R, C). We show the results in Figure 5, where we display the average metric evaluations for each bond length α . In the plot legend, we report the last loss function value during training (complete training loss curves can be found in Appendix B).

From Figure 5, it can be seen that the lowest value of the KL divergence occurs for $\alpha = 7$: this result motivates the usage of such value in Section VB, and suggests that the training is most accurate for this choice of the hyperparameter value. Since the lowest value of the loss function corresponds also to the best validity-based generalization performance, as shown in Figure 5, we notice that we are not reaching the overfitting regime. If this were the case, most likely the model would exhibit a data-copying behaviour, which would imply a lower value of the generalization metrics. We expect that our metrics will be able to identify overfitting behaviours when associated to an extremely successful training curve (Table IV).

As the bond dimension grows, we see an increase in (F, R, C) up to $\alpha = 7$, and then the metrics’ values begin to decrease; at the same time, we observe an opposite trend

in the cost function values. Thus, it seems that we are hitting a trainability *Goldilocks region* around $\alpha \approx 7$, with $\alpha < 7$ leading to underperforming models and $\alpha > 7$ being too expressive for the model to be able to train successfully. These results demonstrate that we can use our metrics to identify thresholds in hyperparameter tuning and to get insights on the trainability of the model, thus potentially spotting training pitfalls associated to its expressivity.

2. Mode collapse in GAN

One of the major issues that affects GAN training is the so-called mode collapse behaviour [53]. This undesired phenomenon occurs when the generator learns to produce a very limited number (sometimes only one) of highly plausible outputs, thus affecting the ability of the generative model to further explore the solution space. Since mode collapse is a well-known pitfall, several strategies have been proposed to mitigate this issue in the context of GANs, among which a promising algorithm is the Wasserstein GAN [58, 59].

We propose an example of how our metrics are able to detect mode collapse, when it occurs. We fine tune our hyperparameters such that the GAN exhibits mode collapse behaviour (see details in Table III in Appendix B) for a fixed training dataset. We run a typical⁵ training of this GAN-MC architecture, and then sample 15 query batches from the trained model to compute our generalization metrics (F, R, C).

We display the validity-based metrics for the GAN and GAN-MC in Table I. For the GAN-MC, we see that fidelity and rate are the ideal value of 1, thus suggesting that the model generates exclusively unseen samples with the desired cardinality. However, the coverage value is close to 0, thus it is far from its ideal threshold, since the model is only able to produce one single pattern and does not have the ability to explore the solution space and cover it as much as possible. Such anomaly in the validity-based generalization metrics’ values is not present if the training of a GAN doesn’t exhibit training pitfalls, as displayed by the GAN results in the same table.

We note that these metrics’ values only capture mode collapse behaviour for models that collapse onto an unseen and valid bitstring. If the model were to collapse onto a seen bitstring (in-training mode collapse), F would be not well-defined and both C and R would equal zero. These metrics’ values would be indistinguishable from the perfect memorization regime. In order to avoid this, one should also compare the number of individual queries generated, $|d_{\text{gen}}|$, to the size of the training set T . This would provide the additional information necessary to detect any form of mode collapse. Expected metrics’ values for various mode collapse behaviour along with other model training pitfalls are displayed in Table

IV in Appendix C. Therefore, our metrics reflect mode collapse upon occurrence and therefore they can provide insights on the training progress of generative models.

In order to better visualize the difference between the two models and detect the mode collapse phenomenon, in Figure 6a we display the cardinality distribution of the generated queries for the two GAN variants under examination: for GAN, the distribution is centered around the correct cardinality but shows a larger spread as compared to the case of GAN-MC, where all the queries satisfy the cardinality constraint. Nevertheless, Figure 6b allows one to identify the occurrence of mode collapse onto an unseen and valid bitstring: the GAN-MC model generates always the same query, as opposed to the diversity of samples retrieved from GAN.

These results demonstrate that we can use our metrics to identify the occurrence of a very well-known pitfall affecting the learning process of GANs, thus providing an insightful tool for the challenging task of monitoring the training of generative models.

D. Evaluating and Comparing Models

We use our quantitative metric-based approach to evaluate the validity-based and quality-based generalization capabilities across different generative models and compare their performance.

We run 30 independent trainings for a fixed training dataset and choose the best run, which we define as the run with the lowest loss function at the end of the trainings. Then, we generate 15 query batches from such trained model, for each of the generative models under examination. We note that while we use a fixed training dataset to compare models, this evaluation method holds across multiple training datasets that could be selected from a specific problem instance. Indeed, each dataset is characterized by the same asset universe, cardinality, and seen portion ϵ , but different datasets can be built by simply drawing independent bitstring subsets from the support of $P(x)$. We perform this analysis in Appendix E, showing that validity-based and quality-based generalization metrics for 15 different training datasets display similar values, thus showcasing the robustness of the models’ behaviour, and the conclusions shown in this work.

For validity-based generalization, we construct $\mathcal{D}_{\text{Train}}$ by sampling from a $P(x)$ that is uniform over the solution space of cardinality-constrained bitstrings, whereas for quality-based generalization, $\mathcal{D}_{\text{Train}}$ is biased with cost-related information, i.e., from $P_{\text{Train}}^{(b)}(x)$, as in Eq. (6). As stated previously, we use one fixed dataset for our evaluation in Section VD 1 and Section VD 2. Post training, $Q = 10^5$ queries are collected from each model for comparison.

1. Validity-Based Generalization

We first show the validity-based generalization results for each type of model. While we present these results as both an

⁵ A typical training instance is identified among 30 independent trainings as the one whose mode collapse shows the correct cardinality and whose last associated Hausdorff distance during training is the median. We highlight that the training is performed via an adversarial strategy, hence we use the Hausdorff distance only as a figure of merit to monitor the training.

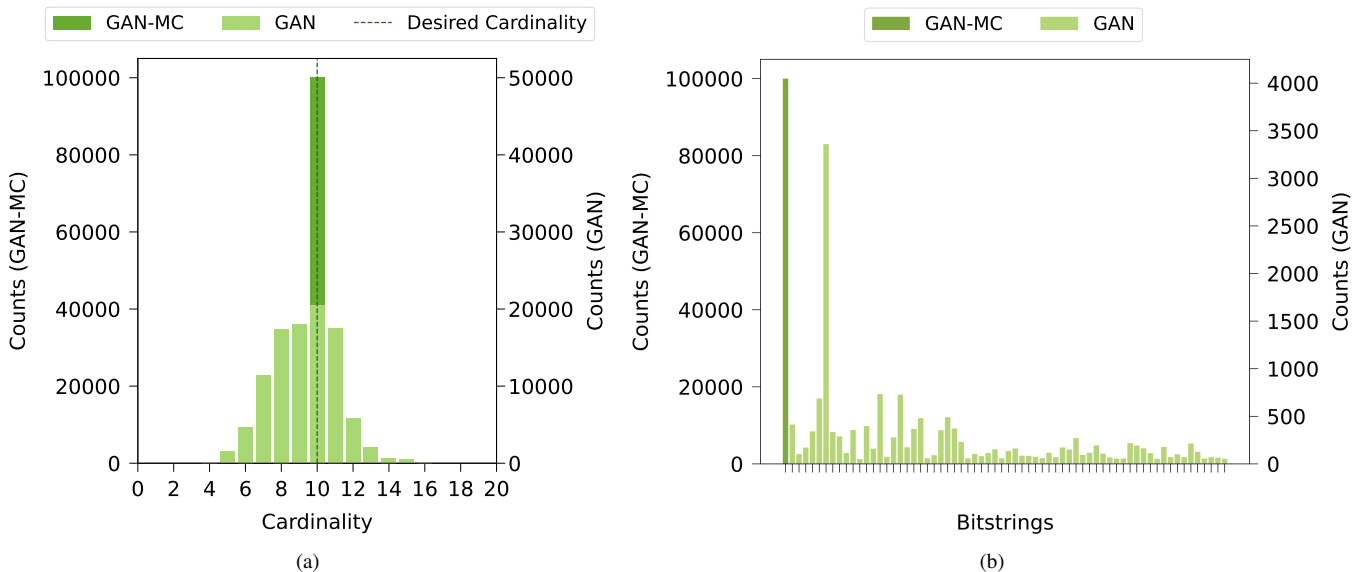


FIG. 6. **Visualization of mode collapse in GAN training.** Figure 6a shows the cardinality distribution of generated queries for GAN and GAN-MC, indicating that GAN-MC produces only samples with the desired cardinality (dashed line), whereas the GAN queries populate a larger subset of the cardinality domain. Hence, GAN-MC is associated to perfect fidelity $F = 1$ and rate $R = 1$. However, in Figure 6b the queries’ diversity is displayed, where the x axis represents the set of distinct generated bitstrings (for readability, bitstrings labels are not shown, and only bitstrings with counts > 50 have been included in the histogram). We can see that GAN-MC always generates the same unseen and valid query (mode collapse phenomenon), as opposed to GAN, which is able to cover a significantly larger portion of the solution space, as reflected by the metrics value in Table I. Note the different scales for the y axes in both Figure 6a and Figure 6b.

evaluation and comparison of models, we would like to emphasize that our results do not speak for all GAN or TNBM models, as each type of model may contain various hyperparameters, multi-layered architectures, and other variances that would lead to different results. Thus, we restrict our comparison to the specific models we trained, as described in Section IV with GAN hyperparameters listed in Table III (Appendix B). We choose to focus on using these models to demonstrate the robustness of our framework and metrics, such that when exploring various GAN, TNBM, or alternative model architectures, this approach can be replicated.

Results for (F, R, C) are listed in Table I, along with the values of the exploration E ; the corresponding results for the metrics’ baseline given by random sampling from the search space are reported in Appendix D. Additionally, we visualize the average validity-based metrics in Figure 7 through a 3D representation. Lastly, Figure 15 in Appendix F gives an intuition of how the two models perform and allows to visualize their different abilities in reconstructing the data distribution $P(x)$, showing the remarkable performance of the TNBM as reflected in the metrics’ values.

In evaluating our models, we see that the TNBM is a clear winner with average values $(0.989, 0.978, 0.409)$. The model achieves near perfect rate and fidelity. As the maximum coverage one can achieve is the number of queries over the size of the solution space ($UB = 0.54$), the TNBM performs remarkably well. Indeed, the ratio of the average coverage to the upper bound UB for the TNBM is high, i.e. $\frac{C}{UB} = 76\%$. However, we note that the upper bound represents a scenario

Metric	TNBM	GAN	GAN-MC	GAN+
E	0.989(0.02%)	0.995(0.02%)	1.0	1.0(0.003%)
F	0.989(0.03%)	0.263(0.6%)	1.0	0.243(0.4%)
R	0.978(0.03%)	0.261(0.6%)	1.0	0.243(0.4%)
C	0.409(0.15%)	0.006(1.7%)	$5.5 \cdot 10^{-6}$	0.001(2.5%)
C/\bar{C}	0.971	0.014	$1.0 \cdot 10^{-5}$	0.002

TABLE I. **Pre-Generalization and validity-based generalization metrics for all models.** We display the average exploration E and the average (F, R, C) values for each best model run with an average and the associated relative percentage error across 15 query batches. All the models exhibit a high exploration rate, thus showing that data-copying is not occurring. We see that our TNBM model outperforms our GAN and GAN+ models by more than 70 percentage points for F and R . C is about 68x larger for TNBM than the GAN models. We further include the ratio of the coverage C to the ideal expected coverage \bar{C} to highlight the large difference between the TNBM and the GAN’s ability to successfully learn the underlying data distribution $P(x)$. Additionally, for GAN-MC, we see perfect F and R and a near zero C value, indicating mode collapse behaviour. Note that no error is provided for the GAN-MC as all the models produce exactly the same values for the metric, except for the coverage whose associated error is negligible.

that would rarely happen, thus representing a pessimistic reference value. A more realistic reference can be derived if one considers the ideal expected coverage \bar{C} when sampling from the data distribution $P(x)$. By means of simple statistical con-

siderations (see e.g., [60, 61]), it can be shown that

$$\bar{C} = 1 - \left(1 - \frac{1}{|S| - T}\right)^Q,$$

and this estimator indicates which coverage C one should expect when the generative model has perfectly learned the data distribution and generates samples accordingly. When comparing the average TNBM coverage to this more realistic reference value, we obtain a surprisingly high value of 97%, which shows that the model has learned an extremely good approximation of the data distribution $P(x)$. In Table I, we include C/\bar{C} values for all models in order to highlight how well each model’s average coverage compares to the ideal expected coverage. The limit of $C/\bar{C} \rightarrow 1.0$ holds for models with perfect generalization.

As shown in Figure 4, the TNBM is able to achieve an improved coverage when sampling up to 3 million queries. The model has a high exploration rate of 98.9%, i.e. $E = 0.989$, such that most of the generated samples were not fed to the model during training. The GAN has much poorer average (F, R, C) values with a slightly higher exploration rate than the TNBM, thus showing that neither of them is performing mere data-copying. The GAN achieves metric values $(0.263, 0.261, 0.006)$, but 99.5% of its generated samples are outside of the training set. One can conclude that while the GAN has the potential to produce novel samples, it requires improved optimization strategies in order to avoid generating noisy samples - i.e. samples that do not match the cardinality constraint - so that fidelity and rate can grow to larger values. The GAN is not able to learn the underlying features as well as the TNBM, and thus is not able to generalize as well. Lastly, we compute the TNBM-to-GAN ratios for the validity-based metrics, and see that the TNBM is $(3.76, 3.75, 68.2) \times$ better than the GAN, respectively across (F, R, C) values. We would like to highlight that using metric ratios, rather than absolute values, allows one to have a clearer picture of the relation between different models, and this strategy is especially useful when considering the coverage, whose absolute value has been shown to be more heavily affected by the number of collected queries Q .

As explained in Section VC 2, we further show visually that our metrics detect mode collapse in GANs. The GAN-MC has an exploration rate of 100% ($E = 1$), demonstrating that the single generated sample was not introduced in the training set. Without the prior knowledge that the model exhibits mode collapse, we can use the average (F, R, C) values $(1.0, 1.0, 5.5e-6)$ to detect this behaviour. If perfect fidelity and rate are achieved, with a coverage near zero, we can conclude that the model has focused in too closely on one or a few unseen and valid bitstrings. In general, whenever $C \rightarrow 0$ we can safely identify the behaviour as mode collapse.

Then, we consider the (F, R, C) values of the GAN+ and see that while the GAN+ is able to explore slightly more than the GAN, the (F, R, C) values are very similar, namely $(0.243, 0.243, 0.001)$, showing that the optimization scheme with Optuna doesn’t bring a significant improvement for our specific GAN model in terms of generalization.

Lastly, we note that F and R are highly correlated for each

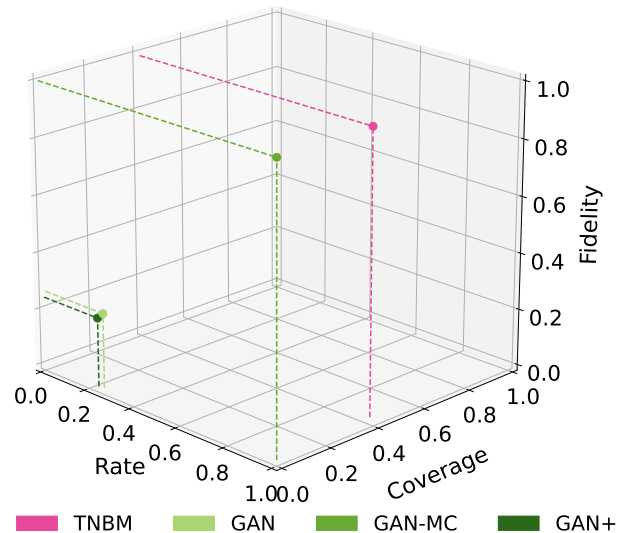


FIG. 7. **3D evaluation of validity-based generalization metrics for different generative models.** The plot displays results for four models, namely the TNBM with $\alpha = 7$ (pink), GAN (light green), GAN-MC (medium green) and GAN+ (dark green). The solid points show the average (F, R, C) values across 15 query batches, whose associated error is too small to be visible in the plot. We see that our TNBM is the clear winner compared to our GAN models.

trained model. This is the case only because in all of the models studied here the exploration E is quite high ($E \approx 1$). In this limit, and given that $R = EF$, then we have $R \approx F$. It is important to note that there is no reason to expect a value of E to be similar across all models, as it happened for the GAN and TNBM explored here.

2. Quality-Based Generalization

We evaluate our generative models’ ability to generate high quality samples using our quality-based approach and metrics. The models (TNBM and GAN) are evaluated across the two *sample quality* metrics described in Section III C: Minimum Value (MV) and Utility (U). Note that for calculating the MV , as discussed in Section III C, five batches of $Q = 10^5$ queries were used. Hence, the total number of query retrievals used to compute this metric is $5 \times$ the number of query sets one would desire for gathering statistics (in our case, $15 \times 5 = 75$ query sets, but this can be adjusted according to the available sampling budget).

When averaged over the 15 independent query retrievals, both the TNBM and the GAN meet the conditions in Eq. (11) and in Eq. (12), as shown in Table II.

We see that our TNBM exhibits a lower MV than our GAN, even though both beat the training set on average. Thus, our TNBM model shows slightly enhanced performance when searching for a minimum value of the cost function C , which is assumed to be the financial risk $\sigma(x)$ in the specific appli-

Metric	TNBM	GAN	Threshold
MV	0.1017(0.01%)	0.1024(0.17%)	0.1035
U	0.1049(0.017%)	0.1048(0.02%)	0.1059

TABLE II. **Quality-based generalization metrics for TNBM and GAN models.** The first column shows values obtained by averaging over all 15 query retrievals for the TNBM’s sample quality performance, along with the associated relative percentage error. The second column displays the metrics’ values and relative percentage error for the GAN model. The last column displays the training threshold, defined as the MV and U computed for the samples in $\mathcal{D}_{\text{Train}}$. We see that both the TNBM and the GAN meet the conditions in Eq. (11) and Eq. (12).

cation we are considering. While this may be relevant when one aims at finding the lowest possible minimum in an optimization task, it may not be the most important condition for alternative tasks that are simply looking for multiple low cost options - not necessarily the lowest. For example, when looking for a large frequency of low cost samples, the condition in Eq. (12) may be more important and robust for comparing models.

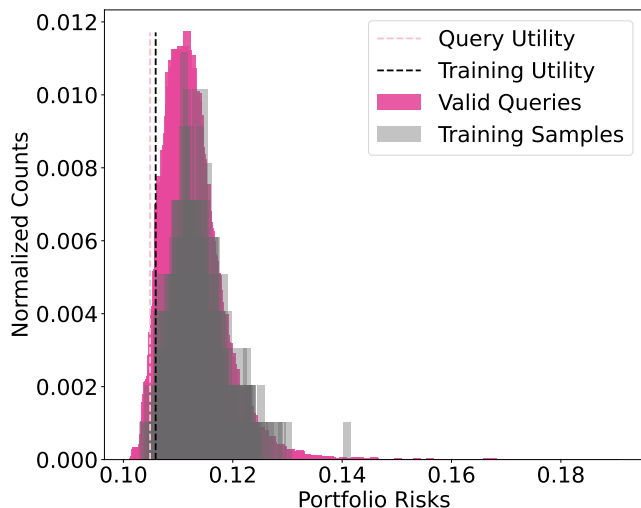


FIG. 8. **Visualization of quality-based metrics for TNBM-generated queries.** The plot displays the number of portfolio counts associated to given risk values. The pink spikes represent valid TNBM queries, whereas the gray spikes represent the samples from the training set. Note that for calculating our metrics, we used $Q = 10^5$ queries, but the training distribution only contains $O(10^3)$ samples. We normalize the counts on the y axis to provide a fair visual comparison between distributions. Because the training distribution is biased towards lower risk values, the model distribution learns this feature in the dataset, and generates an even higher frequency of low risk values. The model queries have a lower utility (pink dashed) than the training set (black dashed), and the model is able to produce samples that have lower risks than those in the training set. We see that our TNBM model is able to effectively generalize to low-risk samples.

From Table II, we observe that the GAN and the TNBM have practically the same U , despite having such a large dif-

ference in (F, R, C) values. We conclude that while both models generate new portfolios that happen to be similarly low in risk when taking the smallest 5% of unseen and valid portfolio risks, the TNBM is simply able to generate more of them than the GAN (TNBM: 4556, GAN: 843, i.e., $5.4\times$). We display these utility samples for TNBM in Figure 8, demonstrating the comparison of U relative to the training distribution P_{Train} . We include the same figure for the GAN in Appendix F.

Hence, the GAN is able to generalize to similarly low risk portfolios as the TNBM, but fewer in number and less diverse than those of the TNBM. Our (F, R, C) metrics support that this generalization diversity is one of the largest differences between our TNBM and our GAN. Therefore, our TNBM model achieves superior performance when looking to produce a large diversified batch of new low-risk valid portfolios. We note that it remains an open question as to why the TNBM’s performance is of such high quality. Investigating the nature of the model’s inductive bias remains an ongoing research effort and opens an interesting opportunity to understand the power of quantum and quantum-inspired model when compared to their classical counterparts.

Lastly, we calculate the number of unique portfolios each model is able to produce that have a lower associated risk than a critical cost in the training set $C'(x)$. When this critical value is equivalent to the sample with the lowest risk in the training set, our TNBM on average is able to beat our GAN with a 61:4 ratio. In other words, our TNBM model is able to generalize to 61 unique portfolios that have a lower risk than the lowest risk in the training set, while the GAN can only produce 4 (i.e., $\sim 15\times$). We introduce this condition in Eq. (13) on top of the other two metrics in order to have an additional layer to determine whether a model is suitable for generalization. Note that one could adjust this critical cost threshold $C'(x)$ to relax the restriction. For example, when $C'(x)$ is equivalent to the risk taken at cutoff of the lowest 5% of samples in the training set, the TNBM-to-GAN ratio becomes 6709:345 on average (i.e., $\sim 19\times$).

While the model might meet the *sample quality* requirements Eq. (11) and Eq. (12), it might be poor at finding many samples with lower cost than $C'(x)$, which is not ideal when one is not only concerned with the global minimum, but also with generating a large quantity of low-cost samples. Our GAN works well under these requirements. On the other hand, our TNBM model shows good quality performance for generalizing to both valid and quality-based portfolios with high diversity and frequency.

SUMMARY AND OUTLOOK

Developing new approaches and frameworks to characterize the generalization capabilities of unsupervised generative models is still a present challenge in both the classical and quantum machine learning community [16, 34–37]. In this work, we contribute to this ongoing effort by first unifying nomenclature for discussing generalization in generative models, next introducing a novel quantitative framework with met-

rics for identifying various generalization behaviours with discrete datasets, and, finally, demonstrating the robustness of our approach by evaluating and comparing the generalization capabilities of two well-known generative models: classical GANs and quantum-inspired TNBMs. We emphasize that to the best of our knowledge, this is the first work that quantitatively compares classical and quantum-inspired models for their generalization capabilities.

Conceptually, we define generalization in an unsupervised generative context as a model’s ability to produce queries that belong to the underlying data distribution $P(x)$, but exist outside of the training distribution $P_{\text{Train}}(x)$. We identify a condition for generalization that a model must meet, known as *pre-generalization*, where the model is able to at least generate samples that are outside the support of $P_{\text{Train}}(x)$.

Furthermore, we introduce a framework for evaluating *validity-based generalization* along with three novel metrics, namely: fidelity F , rate R , and coverage C . Respectively, these metrics provide a 3D evaluation regarding a model’s ability to effectively and efficiently generate, and fully retrieve out-of-training samples that belong to the support of $P(x)$. We highlight that our metrics offer a more concrete picture of a model’s generalization capabilities than that of metrics that do not intentionally reward generalization. We provide further discussion in Appendix A on various evaluation schemes, where most evaluate sample quality and diversity rather than novelty. In contrast, our metrics focus on the quality, diversity, and usefulness of novel samples, narrowing the scope to assessing a model’s generalization capabilities.

Lastly, we introduce sample quality metrics to evaluate *quality-based generalization*, namely Minimum Value (MV) and Utility (U). These metrics provide insight into how well a model can generalize to samples that have associated low cost values (i.e., high quality) for a particular objective function of interest. We demonstrate both our validity-based and quality-based framework on a discrete dataset, associated with a prominent application in finance: cardinality constrained portfolio optimization [20, 21].

After defining our approach, we outline the experiments we performed to demonstrate the robustness of our metrics, and show their ability to spot pitfalls in training, as well as to evaluate and compare models. We see that the relative error of our metrics is consistently small, thus suggesting that our metrics show significant robustness when computed on different sets of queries.

Further enhancing the validity of our approach, we show that our metrics can detect trainability and expressibility issues in TNBM and GAN models. Through varying the bond dimension α of the MPS, we are able to see how the hyperparameters affecting the model’s expressibility impact the model’s ability to generalize. We ultimately see that our metrics can detect trainability issues as we increase the model’s expressibility. Additionally, we see that it is possible to use our metrics to detect trainability issues in GANs, such as mode collapse.

Finally, we demonstrate an evaluation and comparison of our TNBM and GAN models using our metrics. We see that the TNBM is a clear winner with (F, R, C) values approach-

ing the ideal value of 1.0. Our quality-based metrics are able to detect that, while both our GAN and TNBM are able to generalize to better quality portfolios than those in the training set, the GAN generates fewer in number with less diversity than those of the TNBM. We see that a lack of diversity in the generated samples is one of the major drawbacks of our GAN model, and that our TNBM model excels in this generative task. We would like to emphasize that we do not aim to make claims about TNBMs outperforming GANs in absolute terms, as we only considered specific realizations of both models. Instead, these results support and highlight the robustness of our generalization approach, that we hope will be used to move towards more rigorous statements about the performance of different generative models, thus setting an appropriate framework to investigate practical quantum advantage.

In future work, we are looking to use this approach to evaluate and compare the generalization capabilities of alternative models. We see the value in further optimizing the hyperparameters of the GAN architecture, and potentially consider different types of networks such as recurrent neural networks (RNNs) and Variational Autoencoders (VAEs), to push their generalization capabilities. As our framework is tailored towards discrete datasets, we are looking to use this approach in the near future on hybrid and fully quantum generative architectures as well. Previously, it has been a challenge to develop frameworks that can detect generalization in quantum circuits as we are capped with training small-depth circuits [62]. With new meta-learning techniques [63–65] among other pre-training and initialization strategies, one may be able to train larger quantum circuits and use our approach to evaluate generalization. Additionally, demonstrating generalization capabilities on real quantum hardware would open up interesting questions as to how noise may impact the generalization capabilities of the quantum circuit models. Lastly, we can look into further applications where generalization can be assessed and might deliver value.

In summary, besides opening the possibility to quantitatively evaluate the generalization capabilities of generative models, one of the most prominent contributions of this work is the possibility to use this framework to unambiguously define and demonstrate practical quantum advantage in this domain. Generalization is the gold standard for measuring the quality of an ML model. With generative modeling having an edge over supervised ML models in the race for quantum advantage [3], we hope this work opens the possibility to start this race on a solid ground, and on datasets with commercial relevance [48]. As shown here, training GANs and other state-of-the-art classical generative models can be challenging to the point that we report a superior performance from the quantum-inspired generative models used here. Although we expect potentially better results from other classical proposals, there is room as well to improve the quantum-inspired versions explored here. There are also exciting possibilities expected from purely quantum generative models such as Quantum Circuit Born Machines [4], as we will be exploring in future work. We hope this work incites both quantum and classical ML experts to use this framework to enhance the performance and design of their models, in this now quantitative

race towards demonstrating practical quantum advantage in generative modeling.

ACKNOWLEDGMENTS

The authors would like to acknowledge Manuel Rudolph, Vladimir Vargas-Calderón, Brian Dellabetta, and Dmitri Iouchtchenko for providing in-depth manuscript feedback. Additionally, the authors would like to thank Javier Alcazar, Luis Serrano, Luca Thiede, and Riley Hickman for providing ML expertise and advisement. Lastly, the authors would like to thank Chris Ballance for additional project support, and to recognize the Army Research Office (ARO) for providing funding through a QuaCGR PhD Fellowship.

-
- [1] Yann LeCun, Y. Bengio, and Geoffrey Hinton, “Deep Learning,” *Nature* **521**, 436–44 (2015).
- [2] Ian Goodfellow, Yoshua Bengio and Aaron Courville, “Deep learning,” (2016), MIT Press.
- [3] Alejandro Perdomo-Ortiz, Marcello Benedetti, John Realpe-Gómez, and Rupak Biswas, “Opportunities and challenges for quantum-assisted machine learning in near-term quantum computers,” *Quantum Science and Technology* **3**, 030502 (2018).
- [4] Marcello Benedetti, Delfina Garcia-Pintos, Yunseong Nam, and Alejandro Perdomo-Ortiz, “A generative modeling approach for benchmarking and training shallow quantum circuits,” *npj Quantum Information* **5** (2018), 10.1038/s41534-019-0157-8.
- [5] Xun Gao, Eric R. Anschuetz, Sheng-Tao Wang, J. Ignacio Cirac, and Mikhail D. Lukin, “Enhancing generative models via quantum correlations,” (2021), arXiv:2101.08354 [quant-ph].
- [6] Christa Zoufal, Aurélien Lucchi, and Stefan Woerner, “Quantum generative adversarial networks for learning and loading random distributions,” *npj Quantum Information* **5** (2019), 10.1038/s41534-019-0223-2.
- [7] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio, “Generative adversarial networks,” (2014), arXiv:1406.2661 [stat.ML].
- [8] Jie Gui, Zhenan Sun, Yonggang Wen, Dacheng Tao, and Jieping Ye, “A review on generative adversarial networks: Algorithms, theory, and applications,” arXiv:2001.06937 (2020), 10.1109/tkde.2021.3130191.
- [9] Lars Ruthotto and Eldad Haber, “An introduction to deep generative modeling,” *GAMM-Mitteilungen* (2021), 10.1002/gamm.202100008.
- [10] Zhao-Yu Han, Jun Wang, Heng Fan, Lei Wang, and Pan Zhang, “Unsupervised generative modeling using matrix product states,” *PRX* **8**, 031012 (2018).
- [11] Junde Li, Rasit Topaloglu, and Swaroop Ghosh, “Quantum generative models for small molecule drug discovery,” arXiv:2101.03438 (2021), 10.1109/TQE.2021.3104804.
- [12] Nathan Brown, Marco Fiscato, Marwin H.S. Segler, and Alain C. Vaucher, “Guacamol: Benchmarking models for de novo molecular design,” *Journal of Chemical Information and Modeling* **59**, 1096–1108 (2019).
- [13] He Huang, Philip S. Yu, and Changhu Wang, “An introduction to image synthesis with generative adversarial nets,” (2018), arXiv:1803.04469 [cs.CV].
- [14] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila, “Analyzing and improving the image quality of stylegan,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2020) pp. 8110–8119.
- [15] Hao Peng, Christopher S. Gates, Bhaskar Pratim Sarma, Ninghui Li, Yuan Qi, Rahul Potharaju, Cristina Nita-Rotaru, and Ian Molloy, “Using probabilistic generative models for ranking risks of android apps,” *Proceedings of the 2012 ACM conference on Computer and communications security* (2012), 10.1145/2382196.2382224.
- [16] Ahmed M Alaa, Boris van Breugel, Evgeny Saveliev, and Michaela van der Schaar, “How faithful is your synthetic data? sample-level metrics for evaluating and auditing generative models,” arXiv:2102.08921 (2021).
- [17] Eric R. Anschuetz and Cristian Zanoci, “Near-term quantum-classical associative adversarial networks,” *Physical Review A* **100**, 052327 (2019).
- [18] He-Liang Huang, Yuxuan Du, Ming Gong, Youwei Zhao, Yulin Wu, Chaoyue Wang, Shaowei Li, Futian Liang, Jin Lin, Yu Xu, *et al.*, “Experimental quantum generative adversarial networks for image generation,” *Physical Review Applied* **16**, 024051 (2021).
- [19] Manuel S. Rudolph, Ntwali Toussaint Bashige, Amara Katarbwa, Sonika Johr, Borja Peropadre, and Alejandro Perdomo-Ortiz, “Generation of high resolution handwritten digits with an ion-trap quantum computer,” (2020), arXiv:2012.03924 [quant-ph].
- [20] Javier Alcazar and Alejandro Perdomo-Ortiz, “Enhancing combinatorial optimization with quantum generative models,” arXiv:2101.06250 (2021).
- [21] Javier Alcazar, Vicente Leyton-Ortega, and Alejandro Perdomo-Ortiz, “Classical versus quantum models in machine learning: insights from a finance application,” *Machine Learning: Science and Technology* **1**, 035003 (2020).
- [22] Brian Coyle, Maxwell Henderson, Justin Chan Jin Le, Niraj Kumar, Marco Painsi, and Elham Kashefi, “Quantum versus classical generative modelling in finance,” *Quantum Science and Technology* **6**, 024013 (2021).
- [23] Yuxuan Du, Min-Hsiu Hsieh, Tongliang Liu, and Dacheng Tao, “Expressive power of parametrized quantum circuits,” *Physical Review Research* **2**, 033125 (2018).
- [24] Brian Coyle, Daniel Mills, Vincent Danos, and Elham Kashefi, “The born supremacy: quantum advantage and training of an ising born machine,” *npj Quantum Information* **6** (2020).
- [25] Ivan Glasser, Ryan Sweke, Nicola Pancotti, Jens Eisert, and Ignacio Cirac, “Expressive power of tensor-network factorizations for probabilistic modeling,” *Advances in Neural Information Processing Systems* **32**, 1498–1510 (2019).

- [26] Ryan Sweke, Jean-Pierre Seifert, Dominik Hangleiter, and Jens Eisert, “On the quantum versus classical learnability of discrete distributions,” *Quantum* **5**, 417 (2021).
- [27] Shengjia Zhao, Hongyu Ren, Arianna Yuan, Jiaming Song, Noah D. Goodman, and Stefano Ermon, “Bias and generalization in deep generative models: An empirical study,” in *NeurIPS* (2018).
- [28] Matthias C. Caro, Hsin-Yuan Huang, M. Cerezo, Kunal Sharma, Andrew Sornborger, Lukasz Cincio, and Patrick J. Coles, “Generalization in quantum machine learning from few training data,” (2021), arXiv:2111.05292 [quant-ph].
- [29] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals, “Understanding deep learning requires rethinking generalization,” (2017), arXiv:1611.03530 [cs.LG].
- [30] Leonardo Banchi, Jason Pereira, and Stefano Pirandola, “Generalization in quantum machine learning: A quantum information standpoint,” *PRX Quantum* **2** (2021), 10.1103/prxquantum.2.040321.
- [31] David H Wolpert, *The mathematics of generalization* (CRC Press, 2018).
- [32] Amira Abbas, David Sutter, Alessio Figalli, and Stefan Woerner, “Effective dimension of machine learning models,” (2021), arXiv:2112.04807 [cs.LG].
- [33] Ali Borji, “Pros and cons of GAN evaluation measures,” *Computer Vision and Image Understanding* **179**, 41–65 (2019).
- [34] Ali Borji, “Pros and cons of gan evaluation measures: New developments,” arXiv:2103.09396 (2021).
- [35] Mehdi S. M. Sajjadi, Olivier Bachem, Mario Lucic, Olivier Bousquet, and Sylvain Gelly, “Assessing generative models via precision and recall,” in *NeurIPS* (2018).
- [36] Tuomas Kynkäänniemi, Tero Karras, Samuli Laine, Jaakko Lehtinen, and Timo Aila, “Improved precision and recall metric for assessing generative models,” arXiv:1904.06991 (2019).
- [37] Muhammad Ferjad Naeem, Seong Joon Oh, Youngjung Uh, Yunje Choi, and Jaejun Yoo, “Reliable fidelity and diversity metrics for generative models,” in *International Conference on Machine Learning* (PMLR, 2020) pp. 7176–7185.
- [38] Casey Meehan, Kamalika Chaudhuri, and Sanjoy Dasgupta, “A non-parametric test to detect data-copying in generative models,” arXiv:2004.05675 (2020).
- [39] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen, “Improved techniques for training GANs,” *Advances in neural information processing systems* **29** (2016).
- [40] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter, “Gans trained by a two time-scale update rule converge to a local nash equilibrium,” (2018), arXiv:1706.08500 [cs.LG].
- [41] Loïc Simon, Ryan Webster, and Julien Rabin, “Revisiting precision and recall definition for generative model evaluation,” (2019), arXiv:1905.05441 [cs.LG].
- [42] Thomas Hellström, Virginia Dignum, and Suna Bensch, “Bias in machine learning—what is it good for?” arXiv:2004.00686 (2020).
- [43] Tai-Danae Bradley, E M Stoudenmire, and John Terilla, “Modeling sequences with quantum states: a look under the hood,” *Machine Learning: Science and Technology* **1**, 035008 (2020).
- [44] James Stokes and John Terilla, “Probabilistic modeling with matrix product states,” *Entropy* **21** (2019).
- [45] Jacob Miller, Guillaume Rabusseau, and John Terilla, “Tensor networks for probabilistic sequence modeling,” (2020), arXiv:2003.01039 [cs.LG].
- [46] Natalie Wolchover, “New theory cracks open the black box of deep learning,” (2017).
- [47] Marcel Hinsche, Marios Ioannou, Alexander Nietner, Jonas Haferkamp, Yihui Quek, Dominik Hangleiter, Jean-Pierre Seifert, Jens Eisert, and Ryan Sweke, “Learnability of the output distributions of local quantum circuits,” (2021), arXiv:2110.05517 [quant-ph].
- [48] Rubén Ruiz-Torrubiano, Sergio García-Moratilla, and Alberto Suárez, “Optimization problems with cardinality constraints,” in *Computational Intelligence in Optimization: Applications and Implementations*, edited by Yoel Tenne and Chi-Keong Goh (Springer Berlin Heidelberg, Berlin, Heidelberg, 2010) pp. 105–130.
- [49] Qi Gao, Gavin O. Jones, Michihiko Sugawara, Takao Kobayashi, Hiroki Yamashita, Hideaki Kawaguchi, Shu Tanaka, and Naoki Yamamoto, “Quantum-classical computational molecular design of deuterated high-efficiency oled emitters,” (2021), arXiv:2110.14836 [quant-ph].
- [50] David J. C. MacKay, *Information Theory, Inference & Learning Algorithms* (Cambridge University Press, New York, NY, USA, 2002).
- [51] Lucas Theis, Aäron van den Oord, and Matthias Bethge, “A note on the evaluation of generative models,” (2016), arXiv:1511.01844 [stat.ML].
- [52] Harry Markowitz, “Portfolio selection,” *The Journal of Finance* **7**, 77–91 (1952).
- [53] Tong Che, Yanran Li, Athul Jacob, Y. Bengio, and Wenjie Li, “Mode regularized generative adversarial networks,” (2016).
- [54] Kevin Roth, Aurélien Lucchi, Sebastian Nowozin, and Thomas Hofmann, “Stabilizing training of generative adversarial networks through regularization,” in *NIPS* (2017).
- [55] Martin Arjovsky and Léon Bottou, “Towards principled methods for training generative adversarial networks,” arXiv:1701.04862 (2017).
- [56] Mario Lucic, Karol Kurach, Marcin Michalski, Sylvain Gelly, and Olivier Bousquet, “Are gans created equal? a large-scale study,” in *Advances in neural information processing systems* (2018) pp. 700–709.
- [57] Takuya Akiba, Shotaro Sano, Toshihiko Yanase, Takeru Ohta, and Masanori Koyama, “Optuna: A next-generation hyperparameter optimization framework,” in *Proceedings of the 25rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (2019).
- [58] Martin Arjovsky, Soumith Chintala, and Léon Bottou, “Wasserstein generative adversarial networks,” in *International conference on machine learning* (PMLR, 2017) pp. 214–223.
- [59] Ishaan Gulrajani, Faruk Ahmed, Martín Arjovsky, Vincent Dumoulin, and Aaron C. Courville, “Improved training of wasserstein gans,” in *NIPS* (2017).
- [60] Mathematics Stack Exchange, “Expected coverage after sampling with replacement k times,” (2017).
- [61] Mathematics Stack Exchange, “Probability distribution of coverage of a set after x independently, randomly selected members of the set,” (2018).
- [62] Daiwei Zhu, Norbert M Linke, Marcello Benedetti, Kevin A Landsman, Nhung H Nguyen, C Huerta Alderete, Alejandro Perdomo-Ortiz, Nathan Korda, A Garfoot, Charles Brecque, et al., “Training of quantum circuits on a hybrid quantum computer,” *Science advances* **5** (2019), 10.1126/sciadv.aaw9918.
- [63] Guillaume Verdon, Michael Broughton, Jarrod R. McClean, Kevin J. Sung, Ryan Babbush, Zhang Jiang, Hartmut Neven, and Masoud Mohseni, “Learning to learn with quantum neural networks via classical neural networks,” (2019), arXiv:1907.05415 [quant-ph].
- [64] Max Wilson, Rachel Stromswold, Filip Wudarski, Stuart Hadfield, Norm M. Tubman, and Eleanor G. Rieffel, “Optimizing

- quantum heuristics with meta-learning,” *Quantum Machine Intelligence* 3, 13 (2021).
- [65] Frederic Sauvage, Sukin Sim, Alexander A. Kunitsa, William A. Simon, Marta Mauri, and Alejandro Perdomo-Ortiz, “Flip: A flexible initializer for arbitrarily-sized parametrized quantum circuits,” (2021), arxiv:2103.08572 [quant-ph].
- [66] Mikołaj Bińkowski, Danica J. Sutherland, Michael Arbel, and Arthur Gretton, “Demystifying mmd gans,” (2021), arXiv:1801.01401 [stat.ML].
- [67] Ishaan Gulrajani, Colin Raffel, and Luke Metz, “Towards GAN benchmarks which require generalization,” arXiv:2001.03653 (2020).
- [68] Qiantong Xu, Gao Huang, Yang Yuan, Chuan Guo, Yu Sun, Felix Wu, and Kilian Weinberger, “An empirical study on evaluation metrics of generative adversarial networks,” (2018), arXiv:1806.07755 [cs.LG].
- [69] Ryan Webster, Julien Rabin, Loïc Simon, and Frédéric Jurie, “Detecting overfitting of deep generative networks via latent recovery,” in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (2019) pp. 11265–11274.
- [70] Hoang Thanh-Tung and Truyen Tran, “Toward a generalization metric for deep generative models,” (2021), arXiv:2011.00754 [cs.LG].
- [71] Jinchun Xuan, Yunchang Yang, Ze Yang, Di He, and Liwei Wang, “On the anomalous generalization of gans,” (2019), arXiv:1909.12638 [cs.LG].
- [72] Diederik P Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” arXiv:1412.6980 (2014).

Appendix A: Related Works

Generative models are powerful and widespread algorithms, but the evaluation of their performance is an open challenge. A huge variety of metrics have been proposed to evaluate generative models: all the methods attempt to quantify the mismatch between the training and the generated distribution, even though they leverage different approaches. In the following, we give a brief overview of the main strategies, pointing to [33, 34] for a thorough review.

A common approach to evaluate generative models uses statistical divergences, such as the Kullback-Leibler divergence and the Wasserstein distance. Unfortunately, the sample complexity of such quantities scales poorly with the dimensionality of the distribution under examination, proving them inadequate in high-dimensional spaces. To overcome this limitation, alternative evaluation metrics with polynomial sample complexity have been proposed, such as Inception Score (IS) [39], Frechét Inception Distance (FID) [40] and Kernel Inception Distance (KID) [66]. Utilizing neural networks to estimate statistical divergences is another available strategy [67].

The main limitation affecting divergence-based metrics lies in that a single number summary is used to score a model, thus being unable to distinguish its different modes of failure. In light of this consideration, [35] introduced precision and recall as metrics to evaluate generative models, hence proposing a 2D evaluation to disentangle the various scenarios that can arise after training. Follow-up contributions have attempted to extend this idea from discrete to arbitrary probability distribu-

tions [41], and to improve precision and recall definitions and computation [36, 37].

This plethora of methods suggests how challenging it is to evaluate generative models. Evaluate the evaluation metrics themselves is an even more complicated task, despite the paramount importance of choosing the right metric for drawing the right conclusions [51]. [68] addresses such a problem, identifying few necessary conditions that a metric should satisfy in order to qualify as a good performance estimator.

One of these conditions is the ability of a metric to detect overfitting. As highlighted by [69], overfitting is basically equivalent to memorization, i.e. anti-generalization, and it’s not always well defined, despite its importance. While being well established in the context of image classification, notions of generalization are less standardized for generative models. Initial studies on this topic in the context of generative models can be found in [38, 67]. Nonetheless, none of the above metrics is specifically tailored to assessing generalization capabilities, or, in other words, to detect overfitting upon occurrence [70]. So far, very few contributions have been proposed to address the interesting problem of studying and quantifying generalization for generative models. Our work aims at filling this gap: we propose a well-defined approach to generalization, deepening insights gathered from [27], and adequate metrics to quantify such capability, following up on the authenticity metric proposed in [16].

[27] proposed a strategy to analyze generalization in generative models, which consists in probing the input-output behaviour of generative models by projecting data onto carefully chosen low dimension feature spaces. By comparing the training and the generated distribution in these spaces, it is possible to assess whether a model can generate out-of-training samples. However, this contribution focuses only on spotting unseen (i.e., non-memorized) samples, without questioning whether these new samples are meaningful data for the task being solved, or useless noise. [71] hints at this limitation, referring to some of the results in [27] as *anomalous generalization* behaviour, where the generated distribution differs significantly from the training distribution. The approach we propose in this work takes off from these two contributions. It goes deeper into the formal definition of generalization, identifying different regimes that allow us to assess if a generative model can generate samples that are new high quality solutions of the problem at hand. Our approach is able to discriminate between anomalous generalization and generalization to valid and good samples. Inspired by the numerosity feature map proposed in [27], we focus our work on discrete probability distributions. This choice allows us to avoid the introduction of complicated embeddings, which are instead required for most of the evaluation metrics proposed so far.

In addition to defining the approach, we introduce several quantifiable measures of the generalization concepts we formalize. [16]’s proposal of the authenticity metric to identify data-copied samples paved the way for our generalization metrics. We share their starting point that precision and recall are independent of generalization capabilities, as the latter is not properly assessed by the former. Additionally, we share their point on the importance of the novelty feature of

the samples generated by a model. The metrics we propose, though, go beyond the authenticity metric in that they aim at equipping the “novelty space” with estimators that quantify important features, i.e. fidelity, rate and coverage of such a unseen space. The focus of our evaluation metrics revolve around the out-of-training generated samples, disregarding the known data.

To better contextualize our metrics with respect to previous works, we highlight that we share the starting point of [35], hence we propose multiple generalization metrics to disentangle different features and modes of failure. Additionally, our metrics satisfy the conditions expressed in [68]: they are able to detect overfitting and mode collapse. The generalization metrics proposed in this work aim at starting a new thread in evaluating generative models, focused on assessing if they are able to generate new valid and valuable data. We anticipate our metrics to be used along with other evaluation metrics that monitor other desired features of the model, to provide a comprehensive assessment of these powerful data generators.

Appendix B: Training Details

Here, we provide additional details on the training process for both the quantum-inspired and the classical model. The TNBM, whose underlying architecture is an MPS, is trained with a DMRG approach [10] with the negative log-likelihood cost function Eq. (16), and the optimization is performed via Stochastic Gradient Descent with learning rate $\eta = 1e-2$. The number of parameters for the worst case in the TNBM is 1864 for our specific model of $\alpha = 7$. As the bond dimensions for each site are adjusted throughout training, we see that the TNBM does not reach the worst case, and instead has a total number of 1152 parameters. The total number of parameters can be calculated by summing over the squared bond dimensions at each site, and multiplying by a factor of 2.

In Figure 9 we show the training curves for TNBM with several values of the bond dimension α , reporting the KL divergence at each training epoch, that complete the data presented in Section VC 1. Once more, we stress that we can detect trainability issues with our metrics that are confirmed by the learning curves trends. However, if we consider models that are successfully trained, we expect that our metrics should be able to detect the overfitting and underfitting regime when varying the hyperparameters (e.g. the bond dimension α which controls the TNBM expressivity).

In the case of the GAN, the architecture is set to be a feed-forward neural network with linear layers. The generator uses a Gaussian prior, ReLU activation function in the hidden layers and sigmoid cost function in the output layer. The discriminator uses Leaky ReLU activation function in all layers, along with a dropout operation before the final layer. The optimization is performed via the Adam algorithm [72]. The values of the hyperparameters are shown in Table III. The number of total parameters in the GAN is the sum of the parameters in the discriminator and the generator. For our specific architecture, the number of parameters is computed in each layer for the discriminator and generator, respectively. For our GAN with

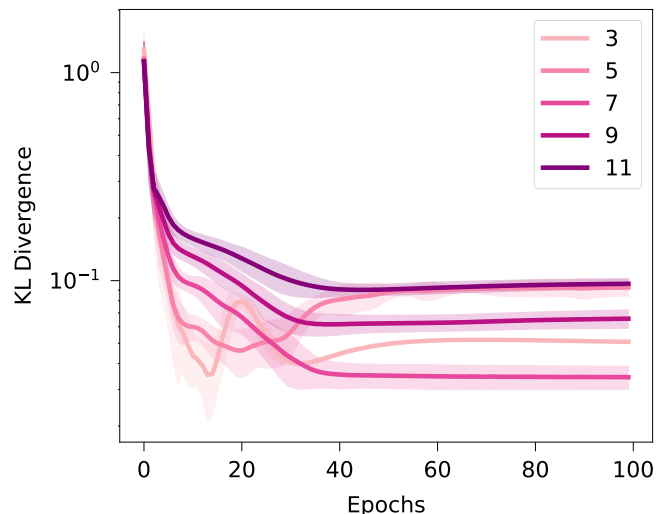


FIG. 9. **TNBM training curves for different bond dimensions.** We plot the KL divergence to monitor the training of the TNBM for bond dimensions $\alpha \in \{3, 5, 7, 9, 11\}$. The lowest KL value is achieved for $\alpha = 7$ after 100 epochs, thus motivating our choice to utilize this value for further studies and model comparisons.

1 hidden layer, we have a total of 4181 parameters.

Hyperparameter	GAN	GAN-MC	GAN+
Prior Size	20	8	12
Hidden Size (G)	20	6	6
Number of Layers (G)	1	4	1
Learning Rate (G)	0.02	0.051	0.001
Hidden Size (D)	20	9	9
Number of Layers (D)	1	3	1
Learning Rate (D)	0.02	0.008	0.006
Negative Slope (D)	0.02	0.007	0.010
Dropout (D)	10^{-5}	0.024	0.107
Batch Size	50	71	56

TABLE III. **GAN hyperparameter values.** The values labelled with $G(D)$ refer to the generator(discriminator). Hidden Size indicates the number of nodes in each hidden layer within G and D , approximated to the same significant digit.

Appendix C: Metrics and Model Behaviours

We provide a short guide to what one could expect to see in our metric values E and (F, R, C) when a model exhibits various training behaviours. This ‘cheat sheet’ can be used to quickly check whether the model is perfectly overfitting/memorizing, perfectly generalizing, exhibiting mode collapse in different nuances, or generating too many novel but noisy samples (i.e., anomalous generalization).

Model Behaviour	E	(F, R, C)	Extra Check
Perfect Generalization	1	(1, 1, 1)	N/A
Perfect Memorization	0	(null, 0, 0)	$ d_{\text{gen}} \sim T$
Anomalous Pre-Generalization	~ 1	(0, 0, 0)	$ d_{\text{gen}} \sim T$
MC (unseen and valid)	~ 1	(1, 1, ~ 0)	N/A
MC (unseen and invalid)	~ 1	(0, 0, 0)	$ d_{\text{gen}} \ll T$
MC (seen and (in)valid)	0	(null, 0, 0)	$ d_{\text{gen}} \ll T$

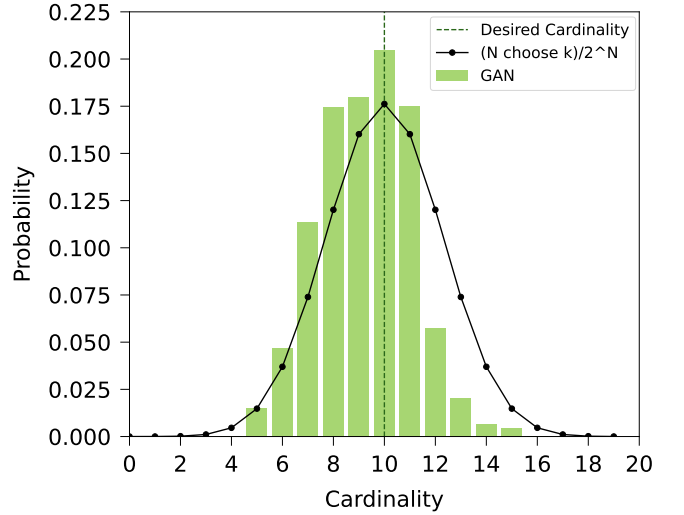
TABLE IV. **Metrics’ values across various model behaviours.** The table displays the E and (F, R, C) values one obtains across different model behaviours such as perfect generalization, perfect memorization/overfitting, generating predominantly noise referred to as anomalous pre-generalization, and mode collapsing (MC) on various bitstring types. We see that F will be null in the cases where the number of unseen generated samples is zero. Additionally, we provide an extra check allowing to distinguish between cases in which the generalization metrics yield the same results.

Appendix D: Random sampling as metrics’ baseline

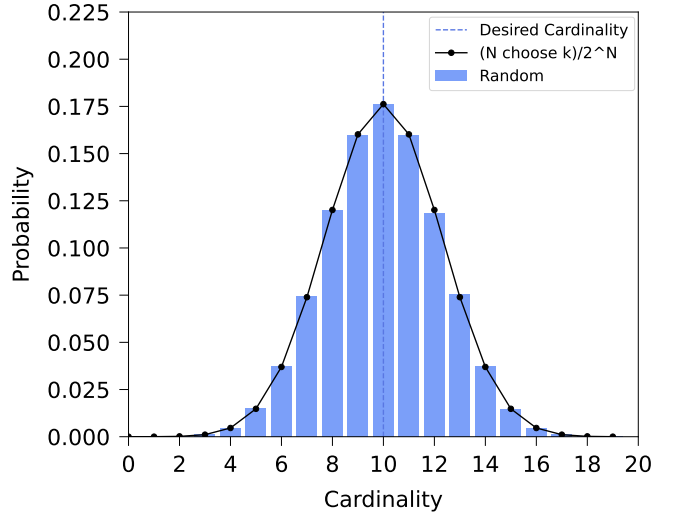
To better characterize the performance of the generative models under examination, we compare their generalization capabilities to a simple baseline: we sample randomly from the search space \mathcal{U} , thus collecting queries to compute the validity metrics, and we compare the results to the ones associated to the TNBM and the GAN. The metrics’ values are summarized in Table V: as expected, both generative models perform better than random sampling, which suggests that during the training process the models were indeed able to learn successfully, despite having different degrees of success. However, the coverage metric in the case of random sampling seems to be higher than the GAN, and this trend persists even considering different numbers of queries Q . What motivates this behaviour is the fact that the GAN suffers from mode collapse: its limited diversity impacts the coverage values, whereas the performance of random sampling is favoured by its higher diversity capabilities. However, Figure 10 shows that the GAN (Figure 10a) is able to generate more samples in the valid space or its vicinity than the random sampler (Figure 10b), thus explaining the higher fidelity of the former as opposed to the latter.

Metric	TNBM	GAN	Random
E	0.989(0.02%)	0.995(0.02%)	0.998(0.013%)
F	0.989(0.03%)	0.263(0.6%)	0.17(0.50%)
R	0.978(0.03%)	0.261(0.6%)	0.17(0.50%)
C	0.409(0.15%)	0.006(1.7%)	0.09(0.48%)

TABLE V. **Pre-generalization and validity-based generalization metrics.** We display the average exploration E and the average (F, R, C) values for each best model run with an average and the associated relative percentage error across 15 query batches. Both the TNBM and the GAN achieve better performance than the random sampler for all the different metrics, except for the GAN coverage as pointed out in the main text.



(a)



(b)

FIG. 10. **Cardinality distribution for GAN and random sampler.** The plots show the percentage of queries with different cardinalities generated by the GAN (Figure 10a) and by the random sampler (Figure 10b). We notice that the GAN is able to produce a higher number of queries with the correct cardinality $k = 10$ (or its vicinity), thus showing that the training process allowed the GAN to partially learn the validity pattern in the training dataset. The black line represents the probability to draw a query with a given cardinality when randomly sampling from the search space \mathcal{U} .

Appendix E: Metrics' Trends

To further demonstrate the power and stability of our metrics, we provide additional details regarding how they scale as we vary the number of queries Q generated from the trained model. Specifically, in Figures 11-14, we plot the values of the validity-based and quality-based generalization metrics and show that most of them do not change with the number of queries - except for coverage, as already shown in Figure 4, and for the minimum value that is displayed in Figure 13. The validity-based trend plots display the constant behaviour of the metrics for both TNBM and GAN as Q increases, along with a dashed black line indicating the ideal metrics value of 1. The quality-based trend plots display the constant behaviour of the utility metric for both TNBM and GAN as Q increases, and a decreasing behaviour for the minimum value as Q increases. The latter is the expected trend: with more queries one has a higher probability of reaching a sample with a lower cost value. For both of these plots, we include a dashed black line indicating the training threshold. This data supports our claim that while our metrics are sample-based, most of them are not dependent on the number of queries.

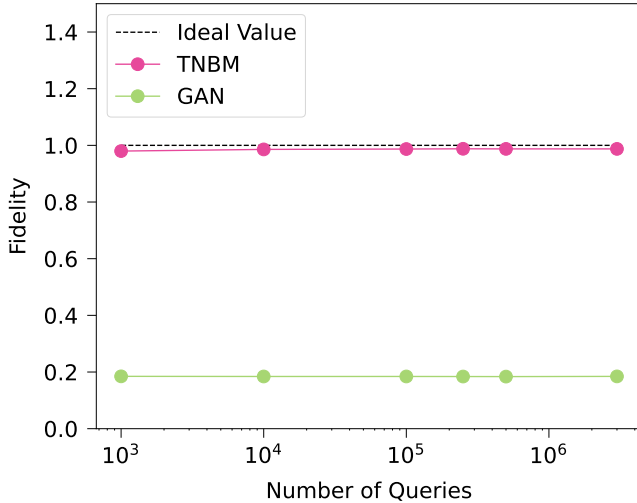


FIG. 11. **Fidelity trends for increasing number of queries.** The plot displays the constant behaviour of the fidelity F for both TNBM (pink) and GAN (green) as we increase the number of queries Q retrieved from the trained models. The dashed black line shows the ideal metric value of 1. In both models, the fidelity is independent of the number of generated queries.

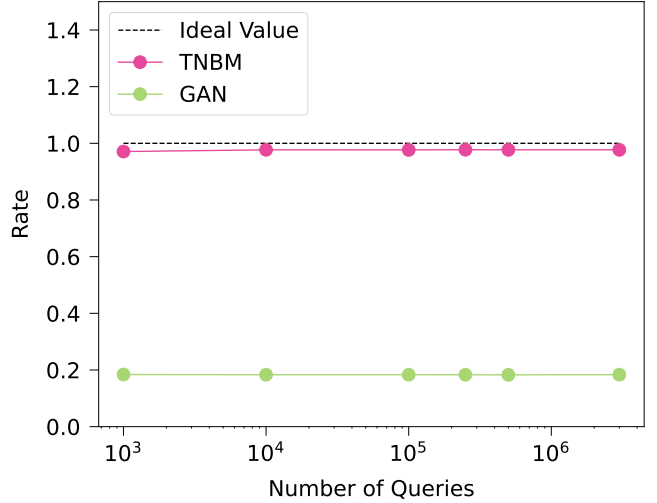


FIG. 12. **Rate trends for increasing number of queries.** The plot displays the constant behaviour of the rate R for both TNBM (pink) and GAN (green) as we increase the number of queries Q retrieved from the trained models. The dashed black line shows the ideal metric value of 1. In both models, the R is independent of the number of generated queries.

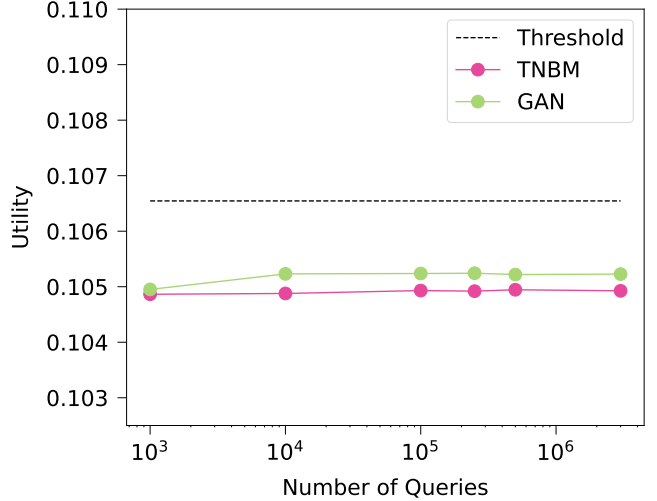


FIG. 14. **Utility trends for increasing number of queries.** The plot displays the constant behaviour of the utility U for both TNBM (pink) and GAN (green) as we increase the number of queries Q retrieved from the trained models. The dashed black line shows the threshold value of the training set U . Both the GAN and TNBM remain under the threshold, independent of the number of queries.

We further propose an investigation on the stability of our approach across various training datasets $\mathcal{D}_{\text{Train}}$. Since a training dataset contains a subset of samples of size T drawn from the solution space \mathcal{S} , it is possible to build different datasets from the same problem instance by randomizing this samples-drawing procedure.

We present the raw data of each of our metrics obtained using 10 distinct datasets built from the same fixed problem

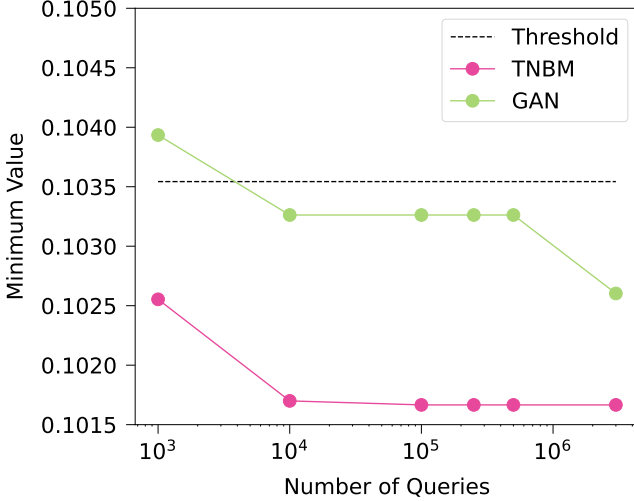


FIG. 13. **Minimum Value (MV) trends for increasing number of queries.** The plot displays the decreasing behaviour of the MV metric for both TNBM (pink) and GAN (green) as we increase the number of queries Q retrieved from the trained models. The dashed black line shows the minimum value of the training set MV . Note that both models not only produce unseen valid samples, but also samples with better quality than those in the training set. As we increase Q , both models are more likely to produce a query with a lower cost value, even if the GAN requires more samples than the TNBM to dip under the threshold.

instance. Thus, all the datasets share the same asset universe, cardinality, and seen portion ϵ as stated in Section V A, and they simply differ for the training bitstrings that get sampled from $P(x)$. Tables VI-E show the results we obtained for the different pre-generalization condition, validity-based and value-based generalization metrics across the 10 different datasets, where each line corresponds to one dataset. We see that the TNBM beats the GAN for all (F, R, C) values. The relative percentage errors across the datasets for (F, R, C) values are smaller for the TNBM: (0.5%, 0.5%, 0.5%) than the GAN: (13%, 13%, 30%), demonstrating that the TNBM produces more stable results across datasets. However, both standard deviations are small enough to show that our metrics produce similar results across various training data.

For the quality-based metrics, we see that the MV for the TNBM is always either equal or less than that of the GAN. However, for U , the TNBM and the GAN trade-off in being the winner. This is not a surprise, as in Table II the TNBM and the GAN produced very similar values for the utility. The same argument from Section V D holds such that both the TNBM and GAN are able to generate low cost samples. Simply, the TNBM contains more diversified high quality samples, which is not captured by the metric U .

E	F	R	C
0.989	0.982	0.971	0.405
0.989	0.978	0.968	0.406
0.989	0.971	0.961	0.401
0.989	0.984	0.973	0.407
0.989	0.983	0.973	0.406
0.989	0.985	0.975	0.407
0.989	0.978	0.967	0.405
0.989	0.977	0.967	0.404
0.989	0.987	0.977	0.406
0.989	0.987	0.977	0.409

TABLE VI. **TNBM pre-generalization and validity-based generalization metrics' values across multiple training datasets from the same problem instance.** We see that the metrics have similar values across the 10 datasets under examination with relative percentage errors (0.5%, 0.5%, 0.5%) for (F, R, C) values respectively. Thus, our metrics produce similar values across multiple training datasets, demonstrating that they are independent of the portion of training samples selected from the valid space.

U	U_T	MV	MV_T
0.1041	0.1064	0.1032	0.1018
0.1040	0.1065	0.1021	0.1034
0.1042	0.1067	0.1019	0.1018
0.1038	0.1064	0.1019	0.1031
0.1029	0.1062	0.1017	0.1033
0.1044	0.1065	0.1018	0.1027
0.1043	0.1068	0.1028	0.1036
0.1048	0.1065	0.1024	0.1029
0.1044	0.1064	0.1017	0.1039
0.1056	0.1064	0.1038	0.1021

TABLE IX. **GAN quality-based metrics' values across various training datasets from the same problem instance.** The second and last columns display the values for the training set, defined as the U and the MV computed for the samples in $\mathcal{D}_{\text{Train}}$. We see that the GAN's U is always less than the training threshold; however, this is not always true for MV , as the GAN has a lower MV value only 70% of the time.

An additional analysis on the stability of the different generative models would be the investigation of their generalization capabilities across different problem instances, especially the ones characterized by larger asset universes, e.g. $N = 500$ (which would correspond to all the assets in the S&P500 index). We highlight here that our approach is not limited to the relatively small universe size considered in this work, i.e. $N = 20$, that was chosen to allow for a practically feasible comparison with quantum generative models in the near term.

E	F	R	C
0.999	0.249	0.249	0.0062
0.996	0.236	0.235	0.0062
0.996	0.309	0.307	0.0063
0.998	0.233	0.233	0.0042
0.995	0.181	0.179	0.0049
0.999	0.232	0.232	0.0061
0.997	0.274	0.274	0.0110
0.997	0.276	0.275	0.0071
0.999	0.239	0.239	0.0066
0.994	0.251	0.249	0.0077

TABLE VII. **GAN pre-generalization and validity-based generalization metrics’ values across multiple training datasets from the same problem instance.** We see that the metrics have similar values across the 10 datasets under examination with mean standard deviations (13%, 13%, 30%) for (F , R , C) values respectively. Despite not being nearly as stable as the TNBM, we see that our metrics produce similar values across multiple training datasets, demonstrating that they independent of the portion of training samples selected from the valid space.

U	U_T	MV	MV_T
0.1049	0.1064	0.1017	0.1018
0.1049	0.1065	0.1017	0.1034
0.1048	0.1067	0.1017	0.1018
0.1049	0.1064	0.1017	0.1031
0.1047	0.1062	0.1017	0.1033
0.1049	0.1065	0.1017	0.1027
0.1051	0.1068	0.1017	0.1036
0.1049	0.1065	0.1017	0.1029
0.1048	0.1064	0.1017	0.1039
0.1049	0.1062	0.1017	0.1021

TABLE VIII. **TNBM quality-based metrics’ values across various training datasets from the same problem instance.** The second and last columns display the values for the training set, defined as the U and the MV computed for the samples in $\mathcal{D}_{\text{Train}}$. We see that the TNBM’s U and MV is always less than the training threshold. Additionally, the same low MV value that exists in the fixed problem universe is generated independent of the training set.

Appendix F: Supplementary Figures

We include supplementary figures to further demonstrate some of our results. Specifically, in Figure 15 we provide 2D visualizations of the data distribution (Figure 15a), the training distribution (Figure 15b), and the output distributions of the trained TNBM (Figure 15c) and GAN (Figure 15d) for a $N = 20$, $k = 10$ problem instance. In the 2D image, every pixel is associated to one of the 2^N bitstrings in the search space \mathcal{U} , and its color encodes the associated probability value. We can see that the bi-dimensional representation of the data distribution displays a non-trivial pattern defined by the solution space \mathcal{S} . Remarkably, provided the small amount of samples that do not demonstrate a very clear pattern in the training distribution, the TNBM and GAN are able to learn the unknown correlations: in particular, the TNBM is able to almost perfectly infer the patterns in the data distribution from very little information.

In Figure 16, we provide a visualization of the GAN quality-based generalization metrics in analogy to Figure 8. By comparing the two plots, we can see that both models reach the low-risk section of the spectrum, but the TNBM samples exhibit more diversity than the GAN ones.

In Figure 17 we display a comparison of the training stability of TNBM and all the three GANs considered in this work, showing how good each of the model is in capturing the correct cardinality pattern encoded in the dataset. We can detect the higher instability affecting GAN models, as opposed to the MPS performance, which appears remarkable. We highlight that even if the TNBM produces only queries with a given cardinality, similar to the GAN-MC histograms, the quantum-inspired model is not exhibiting mode collapse onto an unseen and valid bitstring, as the coverage is not negligible as in the GAN-MC case (see Table I).

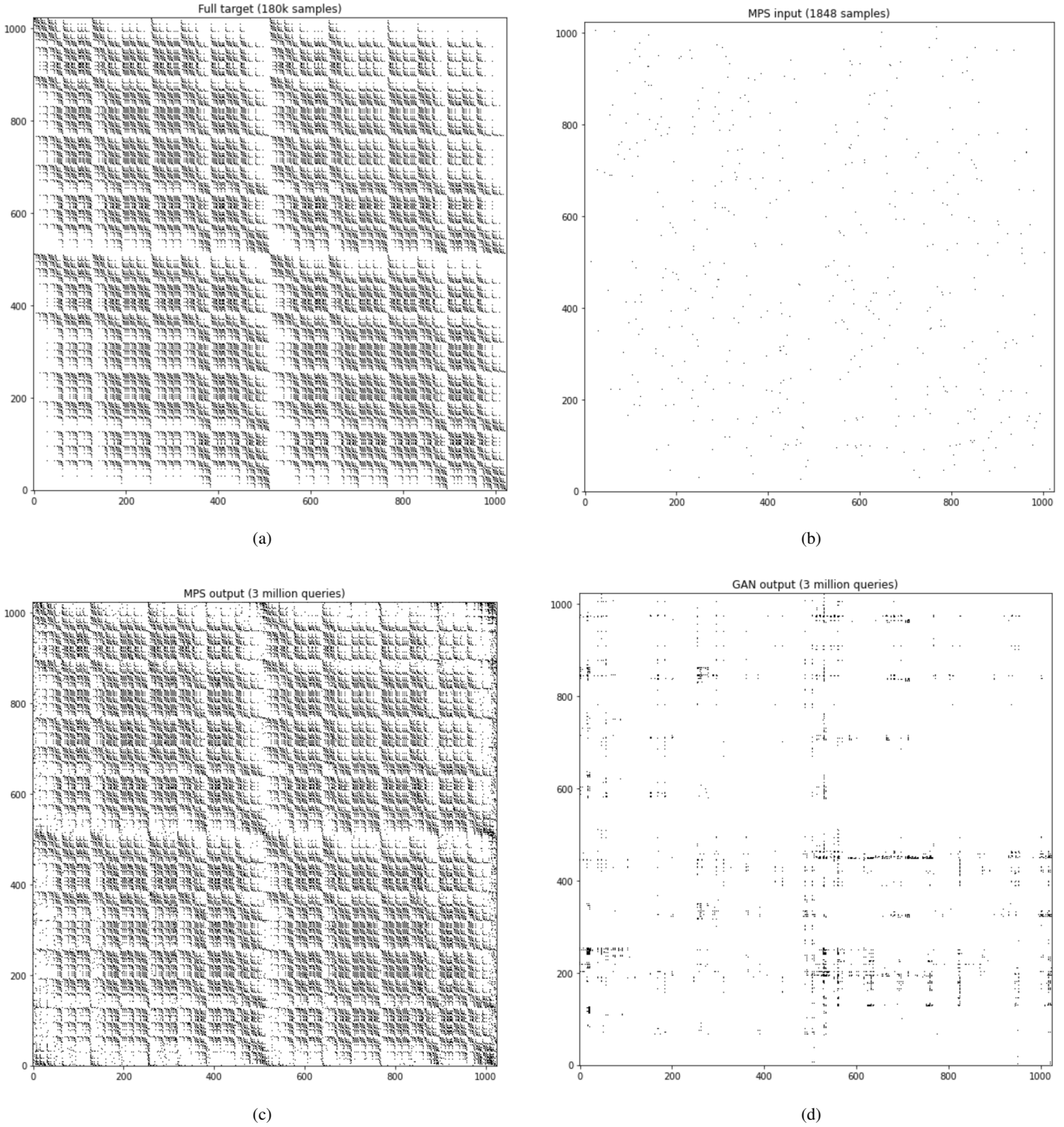


FIG. 15. **2D Visualization of distributions.** Figure 15a shows the 2D visualization of the exact data distribution defined by the solution space S , where we see that a specific pattern emerges from the cardinality. In Figure 15b, we display the 2D visualization for the training distribution, where the same distribution was given to both the TNBM and the GAN models. As shown in Figure 15c, it is very remarkable that with this very limited number of training patterns provided to each model, the TNBM is able to generate the pattern from the data distribution almost exactly (as reflected in the metric values too). On the contrary, in Figure 15d we see that while the GAN is able to learn portions of the pattern, it struggles to reproduce this data distribution.

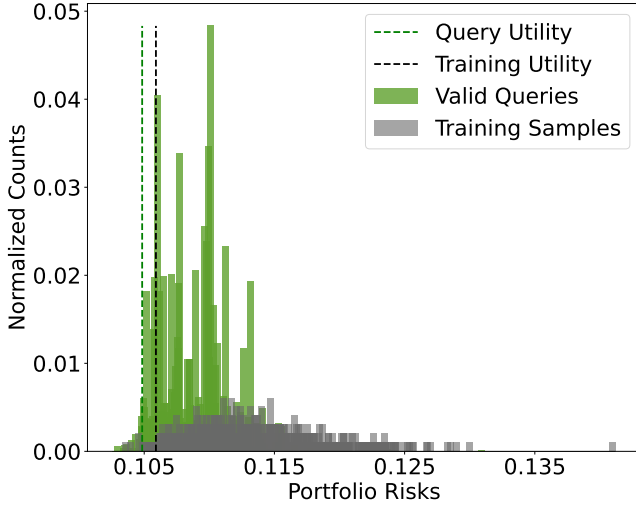


FIG. 16. **Visualization of quality-based metrics for GAN-generated queries.** The plot displays the number of portfolio counts associated to given risk values. The green spikes represent valid GAN queries, whereas the gray spikes represent the samples from the training set. Note that for calculating our metrics, we used $Q = 10^5$ queries, but the training distribution only contains 1,848 samples, hence the need for normalizing the histograms. Similar to the TNBM, the model distribution learns the low-risk bias encoded in the training set, and generates more values of low risk. However, unlike the TNBM, the model frequency counts per query are higher, and the sample diversity is quite low. The queries have a lower utility (green dashed) than the training set (black dashed), thus meeting the condition in Eq. (12). Ultimately, no matter the query count, we see that the GAN can reach low risk queries, but simply has less diversity among them in contrast to the TNBM.

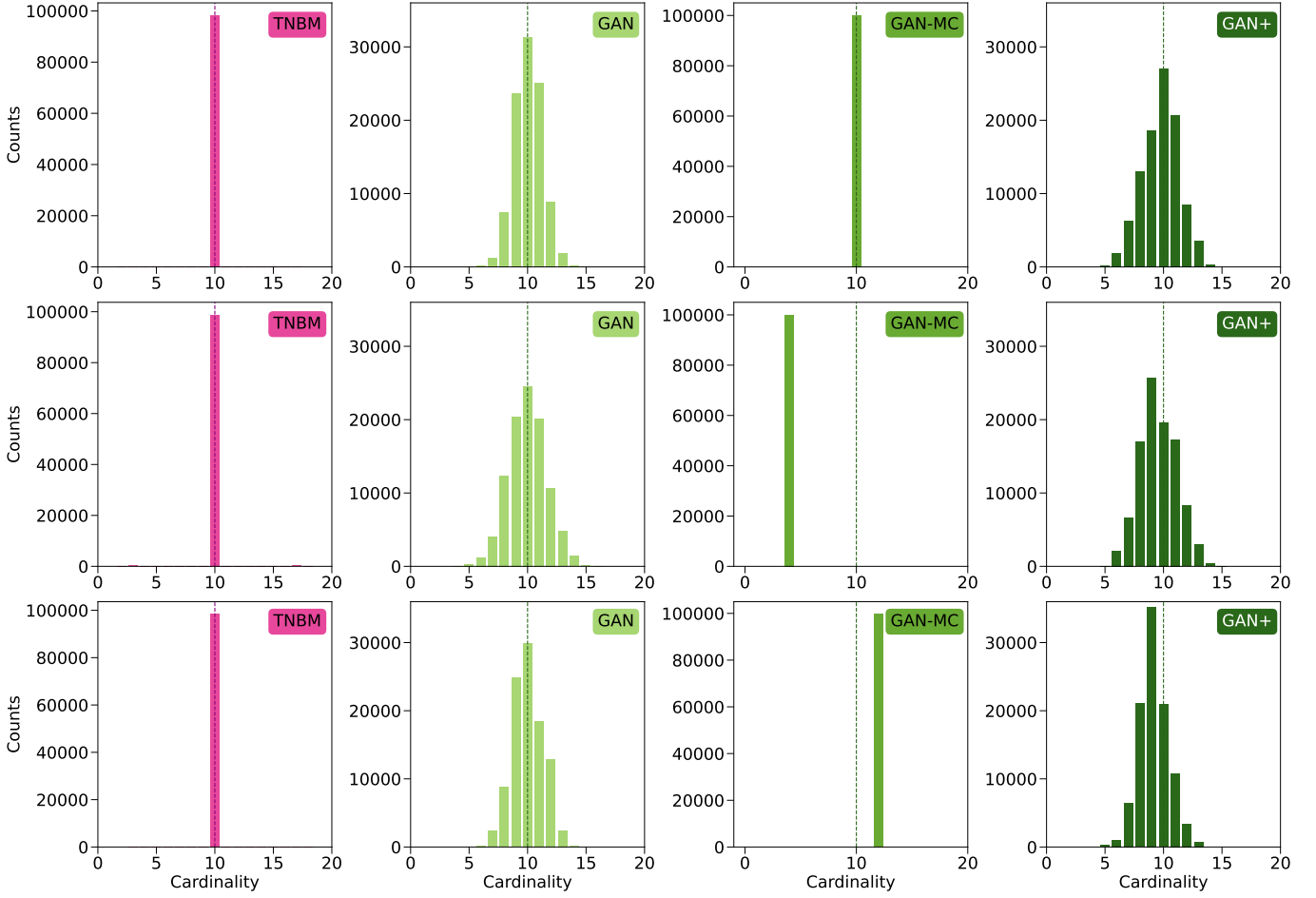


FIG. 17. Cardinality distributions of queries generated by multiple models during independent trainings. We represent cardinality histograms obtained when taking $Q = 10^5$ queries from three independently trained instances of each model family (TNBM, GAN, GAN-MC, GAN+). Each plot displays the cardinality distribution of the retrieved queries, along with the desired cardinality $k = 10$. We can see that the three TNBM models generate queries that always learn accurately the cardinality constraint, whereas the GAN models show less training stability, which is known to be one of the issues affecting this class of classical generative models. Specifically, for GAN and GAN+ we see that while each model always produces at least some valid queries, the centers and tails of the distributions vary greatly for each instance. For GAN-MC, distinct trainings collapse onto different cardinalities, implying that the model is not always guaranteed to generate valid queries.