# cuTN-QSVM: cuTensorNet-accelerated Quantum Support Vector Machine with cuQuantum SDK

Kuan-Cheng Chen**
*QuEST*
*Imperial College London*
London, United Kingdom
kc2816@ic.ac.uk

Tai-Yue Li**
*National Synchrotron Radiation*
*Research Center*
Hsinchu, Taiwan
li.timty@nsrrc.org.tw

Yun-Yuan Wang**
*NVIDIA AI Technology Center*
*NVIDIA Corp.*
Taipei, Taiwan
yunyuanw@nvidia.com

Simon See
*NVIDIA AI Technology Center*
*NVIDIA Corp.*
Santa Clara, CA, USA
ssee@nvidia.com

Chun-Chieh Wang
*National Synchrotron Radiation*
*Research Center*
Hsinchu, Taiwan
wang.jay@nsrrc.org.tw

Robert Wille
*Chair of Design Automation*
*Technical University of Munich*
Munich, Germany
robert.wille@tum.de

Nan-Yow Chen
*National Center for HPC*
*Narlabs*
Hsinchu, Taiwan
nanyow@nchc.narl.org.tw

An-Cheng Yang
*National Center for HPC*
*Narlabs*
Hsinchu, Taiwan
1203087@nchc.narl.org.tw

Chun-Yu Lin
*National Center for HPC*
*Narlabs*
Hsinchu, Taiwan
lincy@nchc.narl.org.tw

*Abstract*—This paper investigates the application of Quantum Support Vector Machines (QSVMs) with an emphasis on the computational advancements enabled by NVIDIA's cuQuantum SDK, especially leveraging the cuTensorNet library. We present a simulation workflow that substantially diminishes computational overhead, as evidenced by our experiments, from exponential to quadratic cost. While state vector simulations become infeasible for qubit counts over 50, our evaluation demonstrates that cuTensorNet speeds up simulations to be completed within seconds on the NVIDIA A100 GPU, even for qubit counts approaching 784. By employing multi-GPU processing with Message Passing Interface (MPI), we document a marked decrease in computation times, effectively demonstrating the strong linear speedup of our approach for increasing data sizes. This enables QSVMs to operate efficiently on High-Performance Computing (HPC) systems, thereby opening a new window for researchers to explore complex quantum algorithms that have not yet been investigated. In accuracy assessments, our QSVM achieves up to 95% on challenging classifications within the MNIST dataset for training sets larger than 100 instances, surpassing the capabilities of classical SVMs. These advancements position cuTensorNet within the cuQuantum SDK as a pivotal tool for scaling quantum machine learning simulations and potentially signpost the seamless integration of such computational strategies as pivotal within the Quantum-HPC ecosystem.

*Index Terms*—Quantum Machine Learning, Quantum Simulation, Tensor Network, cuQuantum SDK, cuTensorNet

**The first three authors contributed equally to this work. The order of these names is based on the alphabetical arrangement of the characters in their names.

## I. INTRODUCTION

In the rapidly evolving landscape of artificial intelligence (AI), machine learning algorithms stand out as pivotal components driving advancements across a multitude of domains [1]. These algorithms, distinguished into supervised and unsupervised learning paradigms, harness the power of data to uncover patterns or make predictions [2]. Supervised learning, in particular, leverages pre-labeled data to train models, with the Support Vector Machine (SVM) being a cornerstone technique in this category [3]. SVMs excel in classifying data into distinct categories by finding an optimal hyperplane in either the original or a higher-dimensional feature space [4]. However, the computational demands of SVMs, especially in the context of large-scale "big data" applications [5], pose significant challenges in terms of both computational resources and execution time.

Enter the realm of quantum computing, a burgeoning field offering profound computational speedups over classical approaches for certain problem types. Among these, Quantum Support Vector Machines (QSVMs) emerge as a promising quantum-enhanced technique for machine learning [6], [7], capable of drastically reducing the computational resources required for SVMs. Leveraging quantum algorithms, QSVMs achieve exponential speedups in both training and classification tasks by performing calculations in parallel and employing quantum-specific optimizations [6], [8]–[10].

However, in the current Noisy Intermediate-Scale Quantum (NISQ) era [11], the practical utility of quantum computers
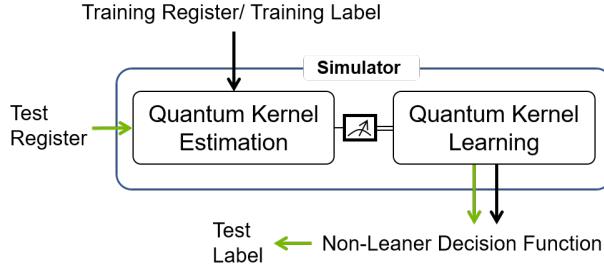
Figure 1. QSVM Simulator: Optimizes quantum kernel estimation and learning, enhancing phase operation and objective evaluation, leading to swift and precise classification outcomes.

is significantly constrained by their availability and imperfect technological state. Challenges such as the fidelity of qubits, the error rates of two-qubit gates, and the limited number of available qubits present substantial hurdles [12]–[14]. Despite the advent of several methodologies aimed at enhancing qubit fidelity—such as Quantum Error Mitigation (QEM) [15], [16] and Dynamical Decoupling (DD) [17], [18]—these limitations persist, impeding the realization of quantum advantage on quantum computing platforms in the current NISQ era [19], [20]. Consequently, the design and validation of quantum-inspired algorithms, or hybrid classical-quantum algorithms, are predominantly conducted through high-performance classical simulations [10], [21]. Furthermore, quantum simulators have shown considerable success in the near-term verification of quantum algorithms on small qubit systems [22].

Within the scope of our research, we have engineered an advanced tensor-network simulation framework, purpose-built to expedite the development of QSVMs through the integration of the cuTensorNet library underlying cuQuantum SDK [23]. This library is meticulously optimized for NVIDIA GPUs and can facilitate QSVM algorithms, requiring noiseless simulations for quantum kernel estimation as depicted in Fig. 1. A pre-computation mechanism is embedded within this workflow, allowing for the reuse of quantum kernel values in the QSVM's complex learning stages, thereby bolstering the efficacy of both the training and classification phases.

Our simulation of the QSVM algorithm, devised for parallel execution using the Message Passing Interface (MPI), harnesses the formidable power of GPU acceleration, equipping our QSVM simulator to efficiently manage vast data sets with only a modest increment in memory requirements. Its flexibility ensures its utility across various quantum machine learning paradigms. Performance benchmarks of our simulator demonstrate that it achieves significant speedups, often exceeding tens of times, which are exponentially better than those achieved by state-of-the-art methods using Qiskit on CPUs [24]. Furthermore, it boasts scalability to multi-GPUs—a testament to the proficiency achieved by the NVIDIA Quantum Team—thereby affirming our simulator's role as a potent and scalable asset in the quantum machine learning arena and the Quantum-HPC ecosystem [21], [23].

This approach's adeptness at handling up to 784 qubits

for QSVMs allows for a comprehensive scaling analysis, shedding light on the quantum acceleration achieved in comparison to traditional classical solvers. These strides in QSVM development signal a major progression towards practical deployment, charting a path for the application of quantum-enhanced methodologies to complex, real-world data classification challenges within the Quantum-HPC Ecosystem [10], [21], [25]–[27]. This paves the way for breakthroughs in quantum computing's applicability and marks a significant contribution to the quantum information sciences.

## II. BACKGROUND

In the research work by Rebentrost et. al. [6], a quantum algorithm is introduced that revolutionizes the computational efficiency of SVMs for big data classification. This QSVM leverages quantum mechanics to achieve logarithmic complexity in terms of both the dimension of the vectors $N$ (qubit number) and the number of training examples $M$ (data size).

Traditionally, SVM is formulated to identify a hyperplane that maximizes the margin between two classes, described in the primal form as:

$$\min_{\mathbf{w},b} \frac{1}{2}||\mathbf{w}||^2 \tag{1}$$

subject to

$$y_j(\mathbf{w} \cdot \mathbf{x}_j + b) \geq 1, \; \forall j, \tag{2}$$

where $\mathbf{w}$ is the normal to the hyperplane, $b$ is the bias, $\mathbf{x}_j$ are the feature vectors, and $y_j$ are the class labels.

The quantum version utilizes the kernel method, mapping the feature space to a higher-dimensional space to allow for nonlinear classification. The kernel matrix $K$, with elements $K_{jk} = k(\mathbf{x}_j, \mathbf{x}_k)$, represents inner products in this space. The quantum algorithm efficiently inverts this kernel matrix, a task that is computationally intensive in classical SVM, especially as $M$ and $N$ grow.

The quantum algorithm for matrix inversion, central to QSVM, significantly reduces the runtime complexity to $O(\log(NM))$ [6]. This quantum routine for non-sparse matrix exponentiation underpins the SVM's optimization, facilitated by the QML algorithm for linear equations, known as the HHL algorithm, which solves $A|x\rangle = |b\rangle$ for Hermitian matrix $A$.

The SVM's dual optimization problem is given by:

$$L(\alpha) = \sum_{j=1}^{M} y_j \alpha_j - \frac{1}{2} \sum_{j,k=1}^{M} \alpha_j K_{jk} \alpha_k, \tag{3}$$

with constraints:

$$\sum_{j=1}^{M} \alpha_j y_j = 0, \tag{4}$$

$$\alpha_j \geq 0, \; \forall j. \tag{5}$$

The QSVM algorithm employs quantum routines to approximate the solution vector $\alpha$, defining the separating hyperplane in the feature space, thus classifying new data points with reduced computational complexity. This approach highlights

the potential of quantum computing to tackle large-scale machine learning problems, characteristic of big data challenges, by utilizing quantum parallelism and efficient manipulation of high-dimensional data.

## III. SIMULATING QSVM USING CUQUANTUM SDK AND CUTENSORNET

### A. Introduction of cuQuantum SDK and cuTensorNet

As the fields of quantum computing and advanced numerical simulations rapidly expand, NVIDIA has introduced cuQuantumK [23], a comprehensive software development kit (SDK), to accelerate quantum circuit simulations with NVIDIA GPUs. The cuQuantum SDK consists of optimized libraries such as cuStateVec and cuTensorNet. cuStateVec is dedicated to state-vector-based simulation methods, enabling substantial acceleration and efficient memory utilization. cuTensorNet, on the other hand, deals with tensor-network-based simulation, where a quantum circuit is represented as a tensor network (Fig.2). In this formulation, the sequence of pairwise contractions plays a role in computational cost as illustrated in Fig.3. The cuTensorNet library provides advanced features for tensor network contractions, including path optimization, approximate tensor network simulation techniques, and multi-GPU multi-node execution. These functionalities facilitate simulations at an unprecedented scale with significant acceleration, empowering researchers to delve into unexplored tensor network theories and complex quantum algorithms. As the demand for more sophisticated computational methods grows within the scientific community, cuStateVec and cuTensorNet stand out for their potential to accelerate research domains across quantum physics, quantum chemistry, and quantum machine learning.
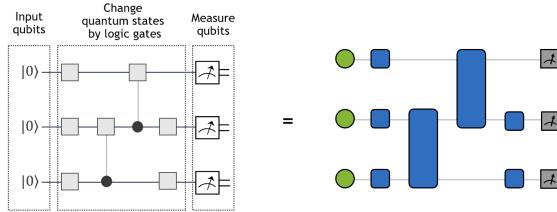


Figure 2. Quantum Circuit Formulation of cuTensorNet

To boost the efficiency of tensor network computation, cuTensorNet delivers modular and finely adjustable APIs, as shown in Fig. 4, tailored for optimizing the pairwise contraction path on the CPU and improving contraction performance on the GPU. This optimization is essential for minimizing both computation cost and memory footprint. The pathfinder workflow is primarily structured in the following manner:

*1) Simplification:* This initial stage focuses on reducing the complexity of the whole tensor network and eliminating redundancies within the network. The implementation involves rank simplification to minimize the number of tensors by removing trivial tensor contractions from the network, resulting in a smaller network for subsequent processing.
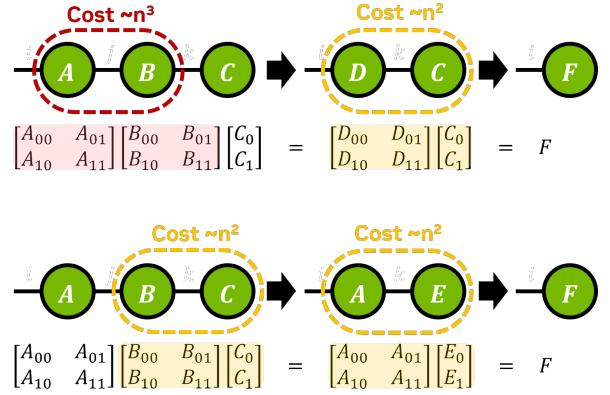


Figure 3. Consider two contraction paths: The upper path results in a higher computational cost

*2) Division:* After simplification, the tensor network undergoes a recursive graph partitioning. This approach segments the network into multiple sub-networks and forms a contraction tree. The binary tree defines the contraction path and can be further optimized at the following stage.

*3) Slicing and Reconfiguration:* The slicing process selects a subset of edges from a tensor network for explicit summation. This technique results in lower memory requirements and allows parallel execution for each sliced contraction. Reconfiguration considers several small subtrees within the full contraction tree and attempts to reduce the contraction cost of the subtrees. cuTensorNet implements dynamic slicing, which interleaves slicing with reconfiguration.

cuTensorNet can efficiently find a contraction path with high quality quickly [23]. This capability accelerates the exploration of quantum mechanics underlying complex systems and quantum machine learning models, particularly in the processing and analysis of high-dimensional data.
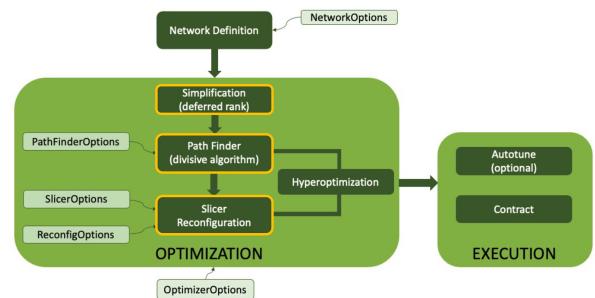


Figure 4. Building blocks of cuTensorNet

Consequently, cuQuantum itself represents a comprehensive suite aimed at bolstering quantum computing simulations. As part of NVIDIA's broader initiative to support the quantum computing ecosystem, cuQuantum provides a range of tools and libraries engineered to optimize various facets of quantum computing simulations, including the programming model CUDA Quantum [23], [28] (CUDA-Q) and frameworks such as Qiskit [29], Pennylane [30], and Cirq. [31]. The

introduction of cuQuantum marks a significant stride towards bridging the theoretical potential of quantum computing with the practical realities of today's computational resources. By offering a scalable and high-performance platform for quantum simulations, cuQuantum not only democratizes access to quantum computing research but also propels the field towards achieving real-world quantum computing applications.

### B. Pipeline of simulating QSVM algorithm

In Fig. 5(a), the depicted pipeline of a QSVM commences with the initial quantum state preparation in a canonical basis state $|0\rangle$. The number of qubits depends on input data features, which can be adjusted using principal components analysis (PCA) to evaluate QSVM with varying qubit counts. A parameterized quantum circuit (QC) follows, designed to map classical data onto a quantum feature space via tunable gates, with the gate parameters refined throughout the training process to facilitate effective data classification. This step, shown in Fig. 5(b), is counterpoised by the introduction of the adjoint of the parameterized quantum circuit ($QC^{\dagger}$), which ensures the reversion of ancillary qubits to their initial state, maintaining the quantum state's coherence for measurement. Central to the QSVM's operation is the quantum entanglement and interference within the high-dimensional feature space, as visualized in the intertwined loops, where the data classification is a result of the constructive and destructive interference patterns that are encoded in the state's amplitude probabilities. The measurement stage yields classical outputs indicative of the data's class membership probabilities, signifying the culmination of the QSVM process—from quantum state preparation and transformation through to interference-based processing and eventual classical information retrieval—for enhanced machine learning efficacy. In this paper, we use a parameterized quantum circuit based on Block-Encoded State (BPS) wavefunctions [32], [33]. This enables QSVM to maintain high classification accuracy even with a greater number of qubits. Worthy of note, our quantum circuit does not decompose into small blocks; instead, we entangle each qubit by adding CNOT gates with linear topology. Such nearest-neighbor connectivity renders the QSVM circuit compatible with near-term quantum hardware.

### C. Complexity of Quantum circuit simulation for QSVM

When executed on classical hardware such as CPUs and GPUs, the simulation of the QSVM algorithm poses significant computational challenges. Figure 6 elucidates these challenges, indicating that the complexity scales exponentially with the number of qubits as $O(2^n)$ and quadratically with data size as $O(n^2)$. Additionally, the complexity scales exponentially with the number of qubits $q$, which encode the feature space within the quantum circuit. This aspect underscores the inherent computational intensity of simulating quantum systems on classical infrastructure.

This scenario highlights the computational complexity advantages that QSVM offers in the realm of quantum machine learning. The simulation demands, in terms of computation
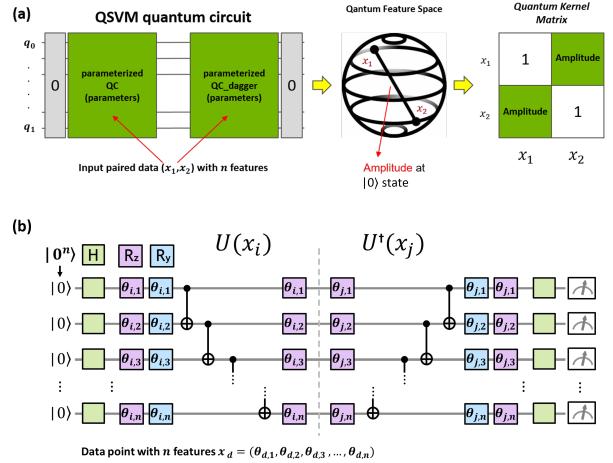


Figure 5. (a) The QSVM pipeline showcasing the quantum circuit transformation of input data into feature space quantum states. (b) A schematic of the QSVM circuit.
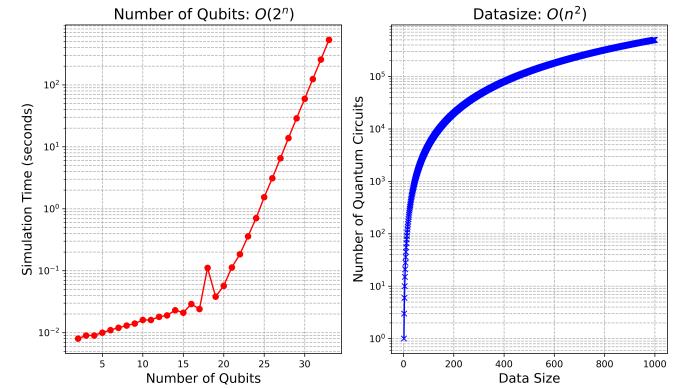


Figure 6. Computational complexity of QSVM simulation. The left graph demonstrates that simulation time scales exponentially with the number of qubits, as $O(2^n)$, while the right graph shows that the number of quantum circuits required scales quadratically with data size, as $O(n^2)$.

time and memory size, grow exponentially with larger datasets and a greater number of qubits, a limitation not encountered when QSVM is run on quantum computers. As demonstrated by Rebentrost et al. [6], the complexity advantage of QSVM can exhibit logarithmic scaling with respect to the product of the number of features and the size of the training set, denoted as $O(\log(NM))$. However, in the NISQ era, the verification of algorithms using traditional CPUs is inevitable. Therefore, this section focuses on leveraging GPU acceleration to address the computational bottlenecks encountered when simulating QSVM with large-scale qubit sizes and processing large datasets.

### D. Simulating QSVM's Kernel Matrix

In this section, we will discuss three methods for simulating a QSVM algorithm as follows. One method involves using a CPU with opt-einsum to optimize the simulation of quantum circuits on CPU devices. The other methods utilize cuStateVector and cuTensorNet for simulating quantum circuits on

**Algorithm 1:** Get Kernel Matrix using opt-einsum

**Input** : Number of data1 $datasize1$, Number of data2 $datasize2$, Circuit einstein summation expression $exp$, List of operands $operands$, Index of data1 and data2 combinations $indices$, Contraction path $path$

1) Initialize $kernel\_matrix \in \mathbb{C}^{datasize1 \times datasize2}$ with all elements set to zero.
2) Set the current operand index $i$ to $-1$.
3) **for** $i_1, i_2 \in \{1, \ldots, indices\}$ **do**
      a) Update the operands index $i \leftarrow i + 1$.
      b) Compute amplitude
         $amp \leftarrow$ opt_einsum.contract($exp, operands[i], path$).
      c) Calculate and store $kernel\_matrix[i_1 - 1][i_2 - 1] \leftarrow \sqrt{\text{amp.real}^2 + \text{amp.imag}^2}$.
   **end**
4) Symmetrize $kernel\_matrix$ by adding its transpose and an identity matrix: $kernel\_matrix \leftarrow kernel\_matrix + kernel\_matrix^T + \text{diag}(\Bbbk_{datasize1})$.

**return** $kernel\_matrix$

---

**Algorithm 2:** Get Kernel Matrix using statevector simulator

**Input** : Number of data1 $datasize1$, Number of data2 $datasize2$, List of quantum circuits $circuits$, Index of data1 and data2 combinations $indices$, statevector simulator $simulator$

1) Initialize $kernel\_matrix \in \mathbb{C}^{datasize1 \times datasize2}$ with all elements set to zero.
2) Set the current operand index $i$ to $-1$.
3) **for** $i_1, i_2 \in \{1, \ldots, indices\}$ **do**
      a) Update the circuits index $i \leftarrow i + 1$.
      b) Save circuits[i] statevector.
      c) Set transpile(circuits[i], simulator).
      d) Run simulator and save result $result$.
      e) Compute amplitude $amp \leftarrow$ result.get_statevector().
      f) Calculate and store $kernel\_matrix[i_1 - 1][i_2 - 1] \leftarrow (\sqrt{\text{amp.real}^2 + \text{amp.imag}^2})$.
   **end**
4) Symmetrize $kernel\_matrix$ by adding its transpose and an identity matrix: $kernel\_matrix \leftarrow kernel\_matrix + kernel\_matrix^T + \text{diag}(\Bbbk_{datasize1})$.

**return** $kernel\_matrix$

---

GPU devices, with acceleration provided by the cuQuantum SDK.

*1) Conventional Simulation of QSVM using CPU/opt-einsum:* Optimized Einsum (opt-einsum) significantly enhances the execution efficiency of einsum-like expressions across various libraries by optimizing contraction orders and employing specialized routines such as BLAS and cuBLAS. Its compatibility with multiple backends, like NumPy, Dask, PyTorch, and TensorFlow, makes it a versatile tool for tensor operations on CPUs.

To ensure a fair comparison between CPU and GPU performance, we utilize the opt-einsum package, which provides optimized tensor computation on CPUs similar to the cuQuantum SDK available for NVIDIA GPUs. The detailed algorithm for simulating the QSVM on CPUs, aimed at equalizing the computational environment to the extent possible, is described in Algorithm 1.

*2) Simulation of QSVM using cuStateVector:* Furthermore, we implement the simulation using the NVIDIA cuQuantum framework, replacing Algorithm 1 with calls to the cuStateVec library to accelerate state vector simulation on GPUs, described in Algorithm 2. The advantage of using cuStateVec includes a speedup of the simulation time by leveraging GPU capabilities and enabling multi-GPU processing with MPI for distributed computing. The effectiveness of cuStateVec in enhancing quantum-circuit-simulation efficiency is evidenced in Lykov et al.'s research work using cuStateVec and the cuQuantum SDK [22].

*3) Simulation of QSVM using cuTensorNet:* However, even with `opt-einsum` facilitating CPU acceleration and `cuStateVector` enabling GPU acceleration, challenges persist due to the complexity of encoding the number of qubits $O(2^n)$ and the size of the data $O(n^2)$. To surmount

---

**Algorithm 3:** Get Kernel Matrix using cuTensorNet with Network Context

**Input** : Number of data1 $datasize1$, Number of data2 $datasize2$, Circuit einstein summation expression $exp$, List of operands $operands$, Index of data1 and data2 combinations $indices$, network options $options$

1) Initialize $kernel\_matrix \in \mathbb{C}^{datasize1 \times datasize2}$ with all elements set to zero.
2) Set the current operand index $i$ to $-1$.
3) Initialize the network with given $options$ to prepare for contraction operations.
4) **for** $i_1, i_2 \in \{1, \ldots, indices\}$ **do**
      a) Update the operand index $i \leftarrow i + 1$.
      b) Reset the network to its initial state before each contraction.
      c) Prepare the operands for contraction based on $i$.
      d) Compute amplitude $amp \leftarrow$ Contract within the network($exp, operands[i], options$).
      e) Calculate and store $kernel\_matrix[i_1 - 1][i_2 - 1] \leftarrow \sqrt{\text{amp.real}^2 + \text{amp.imag}^2}$.
   **end**
5) Symmetrize $kernel\_matrix$ by adding its transpose and an identity matrix: $kernel\_matrix \leftarrow kernel\_matrix + kernel\_matrix^T + \text{diag}(\Bbbk_{datasize1})$.

**return** $kernel\_matrix$

---

these challenges, we present an innovative approach using the `cuTensorNet` library for QSVM simulation. In the creation of the tensor network representation, we seamlessly integrate
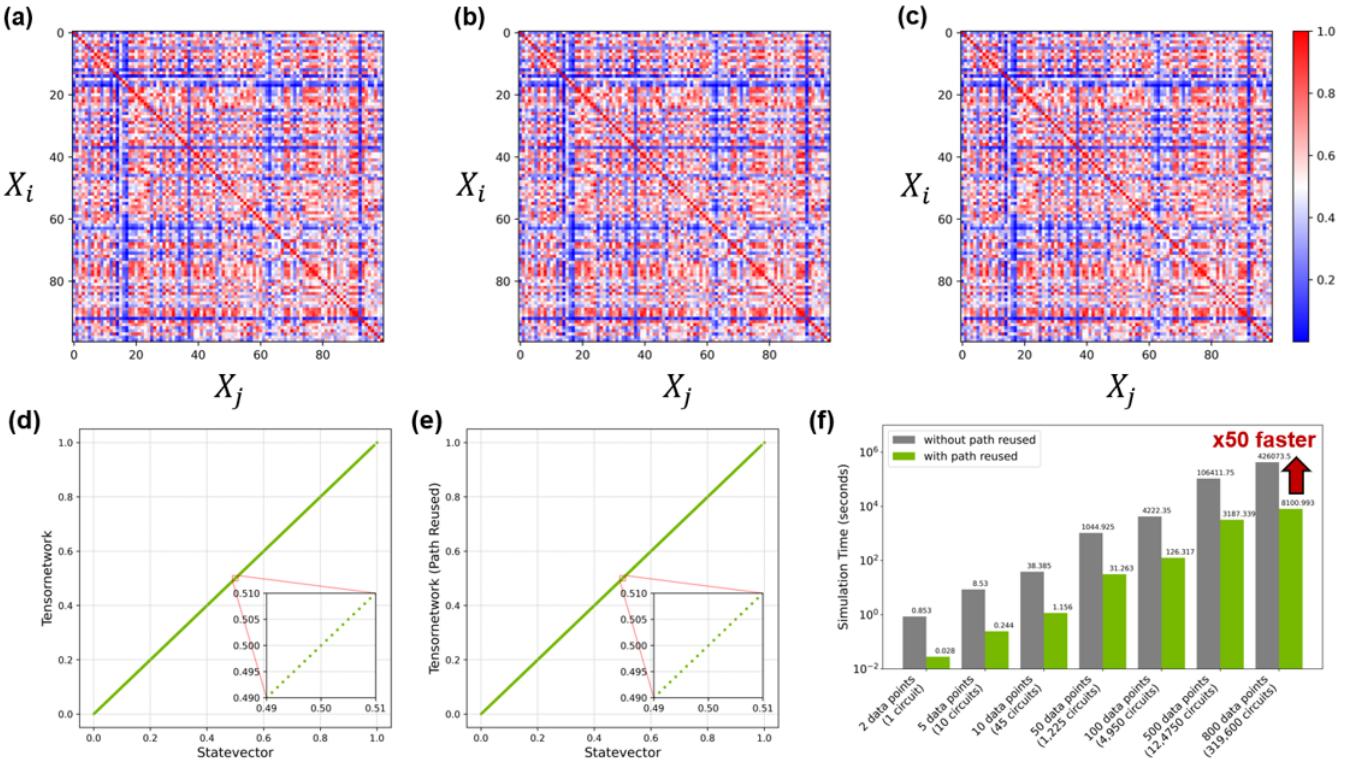
Figure 7. Comparative visualization of quantum kernel matrices and their computation speedups. (a), (b), and (c) illustrate the quantum kernel matrices generated from state vector simulation, tensor network simulation, and tensor network simulation with path reuse strategies, respectively. (d) and (e) feature the parity plots for quantum kernel assessments comparing the outputs of state vector simulations with tensor network and tensor network with path reuse algorithms, demonstrating high concordance. (f) quantifies the performance enhancement attributable to path reuse in tensor network simulations, showcasing significant temporal reductions across an array of dataset sizes.

Qiskit and cuQuantum's built-in `CircuitToEinsum` object. Initially, Qiskit is used to construct a `QuantumKernel` circuit, which is then transformed into 'expression' and 'operand' components by `CircuitToEinsum`. Due to the identical topological structure of the quantum circuit, the same 'expression' component can be reused for subsequent pairs of data. Meanwhile, the 'operand' is updated with parameters from the previously created operand. This approach rapidly transitions data pairs into tensor networks and preserves computational efficiency. The derivation of the kernel matrix—a pivotal component of the SVM—exploits a consistent 'path' to greatly minimize the repetition of contraction order calculations. The detailed algorithm is described in Algorithm 3. This technique not only leverages the computational strength of GPUs but also ensures path reusability, resulting in a considerable acceleration of the simulation process and a dramatic reduction in computational complexity. We will demonstrate those improvements in the next section.

## IV. PERFORMANCE AND BENCHMARKING OF QSVM WITH CUTENSORNET

### A. QSVM Simulation and cuTensorNet-Accelerated QSVM (cuTN-QSVM)

In the outlined simulation workflow, Fig. 1 and 5 illustrate the sequence from the initial input of data to the generation of a quantum circuit for the purpose of encoding. Subsequent steps involve the use of optimized compilation to compute and simulate the quantum circuits, leading to the extraction of a quantum kernel matrix. This matrix is then applied to develop a support vector classifier (SVC).

However, in typical CPU-based workflows, bottlenecks arise in the progression from the construction of quantum circuits to the calculation of the quantum kernel matrix, where the complexity of simulating the QSVM algorithm scales exponentially with the number of qubits, $O(2^n)$, and quadratically with data size, $O(n^2)$. To alleviate these bottlenecks, we incorporate the cuQuantum SDK into QSVM workflow, employing a method of assigned parameters for the formulation of QSVM's quantum circuits. We then maintain a consistent 'expression' for the simulation of these circuits. Ultimately, we apply a 'path reuse' strategy within the tensor network to compute the quantum kernel matrix. This suite of

strategies significantly mitigates the computational complexity associated with processing large datasets, reducing it from $O(n^2)$ to $O(1)$ for pathfinding and substantially enhancing throughput over traditional CPU computations. Importantly, as depicted in Fig. 7, the expressions and paths used in the cuTensorNet during the QSVM simulation process remain unchanged compared to those in CPU and cuStateVector, ensuring that no accuracy is compromised for the sake of expedience. In addition to the path reuse strategy, cuTensorNet offers concurrent execution for tensor network contraction. This technique allows the continued contractions on the GPU asynchronously when tensors are already on the device, thus enhancing computational efficiency by continuing operations without delay. The pronounced speedup achieved through the implementation of path reuse within the cuTensorNet framework is detailed in Fig.7(g), where we report a fiftyfold increase in speed compared to conditions without path reuse.
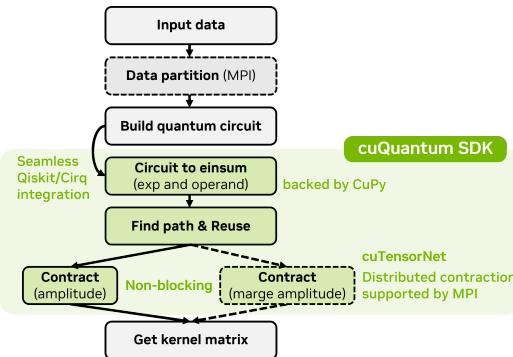


Figure 8. Workflow optimization for QSVM simulation through architectural enhancements, integrating Qiskit/Cirq with cuQuantum SDK. This transition from circuit building to tensor network conversion and kernel matrix computation reduces computational time complexity, leveraging GPU acceleration and multi-GPU strategies.

In the comprehensive workflow outlined in Fig. 8, the input data initiates the construction of a quantum circuit, seamlessly integrated with Qiskit/Cirq. The process advances by converting quantum circuits into tensor networks represented as CuPy arrays, enabling the utilization of in-place operations to efficiently update content for the same operands. Key to enhancing computational efficiency within this framework is the strategic deployment of direct conversion from data to operand, alongside expression reuse for optimizing computational pathways. This step is crucial in minimizing redundancy and ensuring the streamlined execution of the workflow. As the process proceeds, CuPy's capabilities are harnessed to accelerate the computation of the kernel matrix, culminating in the application of the SVC. Moreover, cuTensorNet, as part of the cuQuantum SDK, incorporates advanced strategies such as path reuse and non-blocking operations across multi-GPU configurations.

These approaches significantly reduce the computational overhead from a conventional complexity of $O(2^n)$ to a more scalable $O(n^2)$, thereby enhancing the practicality of executing extensive QSVM simulations with improved pro-

cessing times and efficiency in resource utilization. Fig. 9 illustrates that quantum simulation on the NVIDIA A100 GPU using cuStateVector becomes practically infeasible for more than 50 qubits. However, by employing cuTensorNet, single-contraction simulations can be completed within 0.2 seconds, even with up to 784 qubits. Additionally, Fig. 9 shows that the path reuse strategy can further enhance the speed, offering more than tenfold acceleration when increasing the number of qubits in the QSVM algorithm.

In the GPU-accelerated workflow utilizing cuTensorNet, as delineated in Fig. 8, we are able to expand the feature size (number of qubits) and scale up the data volume for our QSVM algorithm. The evaluation of accuracy resulting from these augmentations will be discussed in the following part, while an in-depth assessment of resource management will be presented in the subsequent section.
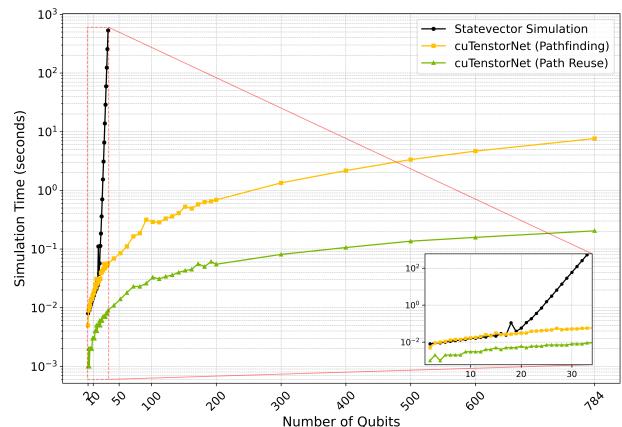


Figure 9. Simulation time comparison for quantum statevector simulation and two cuTensorNet approaches—path plus contraction (no-path reuse), and contraction only (path reuse)—across a range of qubit numbers, up to the equivalent of a 28x28 pixel grid (784 qubits).
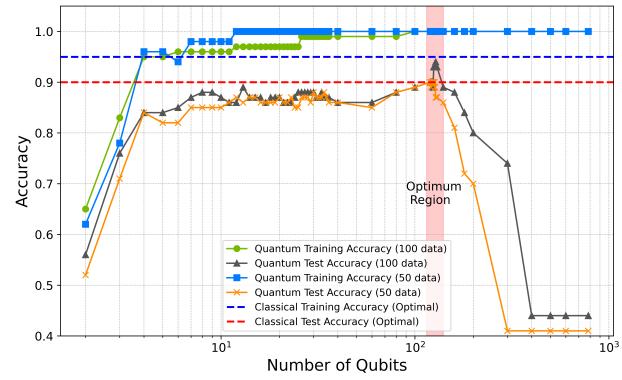


Figure 10. Train and test accuracy of quantum and classical models as a function of the number of qubits. The plot compares quantum model performance with varying training set sizes (50 and 100 data points) against the optimal classical model.

## B. Accuracy Benchmarking

In our study's accuracy benchmarking, we observed that the trained QSVM demonstrates an improvement in accuracy with increasing dataset sizes shown in Fig. 10. Additionally, there tends to be a positive correlation between the number of qubits and accuracy; a larger quantity of qubits typically yields higher accuracy. In the specific case under consideration, the optimized configuration entailed an algorithm using 128 qubits on a sample set of 100. This experiment also indicates that the accuracy sharply declines beyond roughly 200 qubits, attributed to the off-diagonal kernel matrix elements approaching zero. Under optimized parameters, the QSVM achieved test and training accuracies of 94% and 100%, respectively, outperforming the traditional SVM, which achieved accuracies of 90% in testing and 95% in training with a comparable amount of data. These empirical results underscore the enhanced capability of QSVM over traditional SVMs as data volume increases, especially in big data analytics applications, a finding that is corroborated by the work of Chen et al. [10].

Furthermore, in Fig. 11, we extend our evaluation to larger data sizes, focusing on the MNIST dataset, particularly on digits 0 and 9, which are inherently challenging to distinguish. Our analysis reveals that with a limited training dataset size (i.e., 10 instances), the accuracy hovers near 50%. However, as the volume of training data surpasses 100 instances, the accuracy significantly increases to approximately 95%, surpassing the performance of the optimal classical SVM.
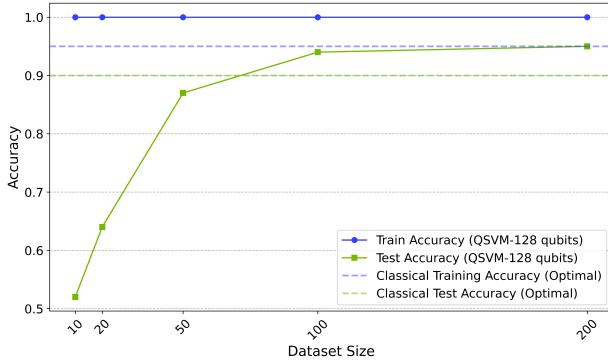


Figure 11. The relationship between dataset size and accuracy for a 128-qubit QSVM algorithm, compared to the optimal classical model.

## C. Simulation with Single CPU and GPU

In this section, we compare the performance of a CPU and a GPU, as illustrated in Fig.12. To ensure a fair comparison, we employed Opt-Einsum for the contraction process on a single AMD EPYC 7J13 CPU, contrasting this with a single NVIDIA A100 GPU using cuTensorNet for the contraction process, with path reuse implemented. The detailed pseudocodes are discussed in Section III-D. Moreover, it was necessary to synchronize the contraction paths in Opt-Einsum with those of cuTensorNet to ensure consistency. As depicted in Fig.12, the speedup provided by the GPU relative to the CPU becomes more pronounced as the number of simulated qubits increases.

Consequently, for large-scale qubit simulations, GPUs demonstrate enhanced scalability and promise substantial benefits for future advanced qubit algorithms in simulation and emulation.
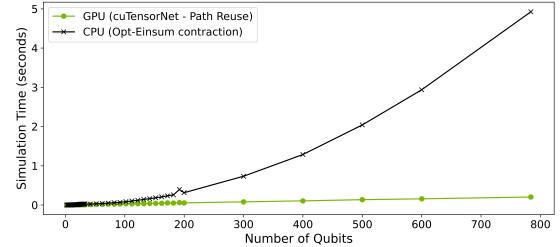


Figure 12. Benchmark QSVM circuit simulation time using a single CPU and a single GPU.

## V. DISTRIBUTED SIMULATION AND RESOURCE ESTIMATION IN HPC

In the final section of our study, a multi-GPU instance was utilized to expand the QSVM model to accommodate a dataset comprising 1,000 data points of 28x28 MNIST images. The implementation of multi-GPU resources to enhance quantum circuit simulation via cuStateVector is thoroughly detailed in the research conducted by Shaydulin et al. [34]. Our emphasis lies on leveraging the data from these experiments to rigorously assess both the computational costs and the temporal demands inherent in simulating the QSVM algorithm within a multi-GPU processing framework.

In our computational environment, each GPU within a node is interconnected using the high-bandwidth NVLink network, optimizing intra-node communication. Inter-node data transfer mandates that information from the GPU be first relayed to the CPU, then onward to the target node, a process that requires careful selection of a communication protocol sensitive to the GPU's positional context. The MPI offers built-in capabilities for such operations, which are enabled by setting the $MPI\_GPU\_SUPPORT\_ENABLED$ environment variable. By harnessing these integrated technological benefits, we have successfully actualized the accelerated computational outcomes for managing large-scale qubit systems and extensive datasets, as illustrated in Fig.13. Comparative analysis indicates that our performance metrics are on par with distributed simulation results documented in the existing scientific corpus, as cited in Bayraktar et. al.'s and Lykov et. al.'s work [22], [23].

## A. Benchmarking cuQuantum Multi-GPU with MPI

Fig. 13 illustrates the execution time required for quantum simulations in relation to the number of qubits. The data compares the performance of a single A100 GPU to systems utilizing 2, 4 and 8 GPUs in conjunction with MPI and within a single NVIDIA DGX node. It is evident from the results that the incorporation of multi-GPUs significantly decreases computation time, highlighting the strong linear speedup of

cuTenserNet with MPI. The trend indicates a substantial reduction in execution time as the number of GPUs is increased, affirming the efficacy of multi-GPU setups in handling large datasets.
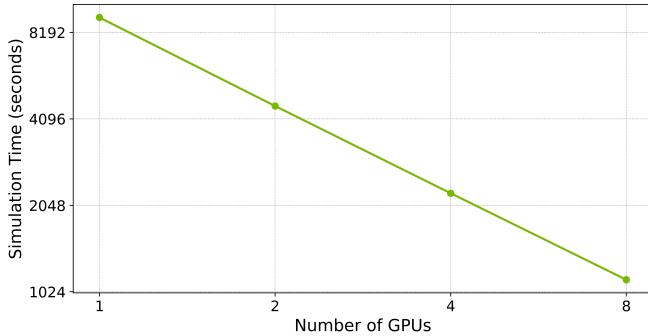


Figure 13. Strong scaling of the QSVM simulation is observed for 1,000 data points across 1, 2, 4, and 8 GPUs, demonstrating linear speedup.

### B. Large Dataset Processing with Multi-GPU instance

Figure 14 presents a comparative analysis of computational time across different configurations, ranging from a single GPU (A100, 80GB) to 2, 4, and 8 multi-GPU arrangements using MPI for processing datasets of various sizes. The results distinctly highlight the superior efficiency and scalability of multi-GPU systems, especially when managing large-scale datasets. A notable reduction in processing time is observed with the integration of an 8-GPU setup, underscoring the considerable advantages of parallel computing for large-scale data analysis. In Figure 14, experimental data (solid line) from 40 to 1,000 data points is extrapolated to estimate the processing time for 10,000 data points, corresponding to nearly 50 million circuits (dashed line). The projection indicates that an eight-GPU system could achieve linear acceleration, reducing a week-long processing task using the simulated QSVM to approximately one day (blue line).
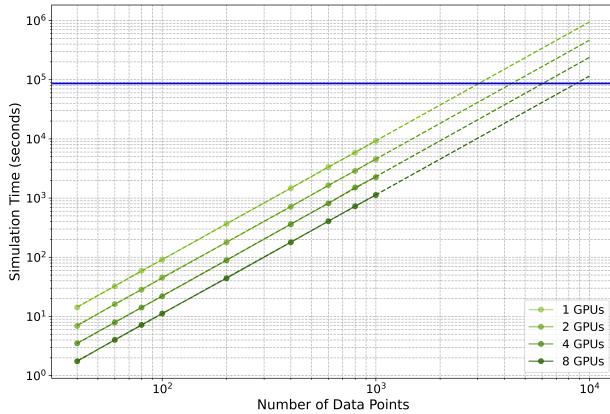


Figure 14. Execution time for quantum simulations against qubit count for a single A100 GPU and MPI-based 2, 4 and 8 multi-GPU setups. The performance enhancement with additional GPUs is evident, underscoring the benefits of parallelized computation.

## VI. Conclusion

This paper has presented a comprehensive study on the application and efficacy of cuTN-QSVMs in conjunction with NVIDIA's cuQuantum SDK. Our findings have consistently shown that by integrating cuTensorNet within our computational workflow, we have been able to significantly decrease the computational complexity involved in simulating QSVM algorithms, especially when working with large qubit counts. Through rigorous performance benchmarking, our approach has demonstrated not only enhanced simulation speeds and efficiency but also scalability across multi-GPU platforms. Additionally, we have observed substantial improvements in accuracy with increasing data set sizes, underlining the potential of quantum computing methods in tackling complex, large-scale data challenges. The integration of cuTensorNet and MPI-based multi-GPU systems reveals a path forward for quantum computing that offers practical advantages for both current research and future applications, bridging the gap between quantum hardware and classical computing resources. This points towards a promising direction for high-performance quantum simulations within the Quantum-HPC ecosystem. Furthermore, we utilized cuTN-QSVM to classify bone cell images captured by synchrotron transmission X-ray microscopy (TXM) at the BL01B1 beamline of the Taiwan Light Source (TLS) in Hsinchu, Taiwan. Bone cell image data possess a high-dimensional feature space but limited data quantity. Traditional machine learning struggles with classifying such data, but QSVM overcomes this, showcasing its advantages. Due to page limit, the classification results of bone cells will be elaborated in the conference.

## Code Visability

The source code for generating the dataset and figures presented in this research work is openly accessible at https://github.com/Tim-Li/cuTN-QSVM.git.

## References

[1] M. I. JORDAN AND T. M. MITCHELL, *Machine learning: Trends, perspectives, and prospects*, Science **349** (2015), pp. 255–260.

[2] D. J. MACKAY, *Information theory, inference and learning algorithms*, Cambridge university press, 2003.

[3] M. A. HEARST, S. T. DUMAIS, E. OSUNA, J. PLATT AND B. SCHOLKOPF, *Support vector machines*, IEEE Intelligent Systems and their applications **13** (1998), pp. 18–28.

[4] B. GHADDAR AND J. NAOUM-SAWAYA, *High dimensional data classification and feature selection using support vector machines*, European Journal of Operational Research **265** (2018), pp. 993–1004.

[5] B. GAYE, D. ZHANG AND A. WULAMU, *Improvement of support vector machine algorithm in big data background*, Mathematical Problems in Engineering **2021** (2021), pp. 1–9.

[6] P. REBENTROST, M. MOHSENI AND S. LLOYD, *Quantum support vector machine for big data classification*, Physical review letters **113** (2014), p. 130503.

[7] X. WANG, Y. DU, Y. LUO AND D. TAO, *Towards understanding the power of quantum kernels in the nisq era*, Quantum **5** (2021), p. 531.

[8] Z. LI, X. LIU, N. XU AND J. DU, *Experimental realization of a quantum support vector machine*, Physical review letters **114** (2015), p. 140504.

[9] C. DING, T.-Y. BAO AND H.-L. HUANG, *Quantum-inspired support vector machine*, IEEE Transactions on Neural Networks and Learning Systems **33** (2021), pp. 7210–7222.

[10] K.-C. CHEN, X. XU, H. MAKHANOV, H.-H. CHUNG AND C.-Y. LIU, *Quantum-enhanced support vector machine for large-scale stellar classification with gpu acceleration*, arXiv preprint arXiv:2311.12328 (2023).

[11] J. PRESKILL, *Quantum computing in the nisq era and beyond*, Quantum **2** (2018), p. 79.

[12] M. KJAERGAARD ET AL., *Superconducting qubits: Current state of play*, Annual Review of Condensed Matter Physics **11** (2020), pp. 369–395.

[13] C. D. BRUZEWICZ, J. CHIAVERINI, R. MCCONNELL AND J. M. SAGE, *Trapped-ion quantum computing: Progress and challenges*, Applied Physics Reviews **6** (2019).

[14] S. J. EVERED ET AL., *High-fidelity parallel entangling gates on a neutral-atom quantum computer*, Nature **622** (2023), pp. 268–272.

[15] Z. CAI ET AL., *Quantum error mitigation*, Reviews of Modern Physics **95** (2023), p. 045005.

[16] K.-C. CHEN, *Short-depth circuits and error mitigation for large-scale ghz-state preparation, and benchmarking on ibm's 127-qubit system*, in 2023 IEEE International Conference on Quantum Computing and Engineering (QCE), vol. 2, IEEE, 2023, pp. 207–210.

[17] A. M. SOUZA, G. A. ALVAREZ AND D. SUTER, *Robust dynamical decoupling for quantum computing and quantum memory*, Physical review letters **106** (2011), p. 240501.

[18] M. PREST AND K.-C. CHEN, *Quantum-error-mitigated detectable byzantine agreement with dynamical decoupling for distributed quantum computing*, arXiv preprint arXiv:2311.03097 (2023).

[19] J. W. Z. LAU, K. H. LIM, H. SHROTRIYA AND L. C. KWEK, *Nisq computing: where are we and where do we go?*, AAPPS bulletin **32** (2022), p. 27.

[20] S. CHEN, J. COTLER, H.-Y. HUANG AND J. LI, *The complexity of nisq*, Nature Communications **14** (2023), p. 6001.

[21] K.-C. CHEN, X. LI, X. XU, Y.-Y. WANG AND C.-Y. LIU, *Quantum-hpc framework with multi-gpu-enabled hybrid quantum-classical workflow: Applications in quantum simulations*, arXiv preprint arXiv:2403.05828 (2024).

[22] D. LYKOV, R. SHAYDULIN, Y. SUN, Y. ALEXEEV AND M. PISTOIA, *Fast simulation of high-depth qaoa circuits*, in Proceedings of the SC'23 Workshops of The International Conference on High Performance Computing, Network, Storage, and Analysis, 2023, pp. 1443–1451.

[23] H. BAYRAKTAR ET AL., *cuquantum sdk: A high-performance library for accelerating quantum science*, in 2023 IEEE International Conference on Quantum Computing and Engineering (QCE), vol. 1, IEEE, 2023, pp. 1050–1061.

[24] D. FORTUNATO, J. CAMPOS AND R. ABREU, *Mutation testing of quantum programs: A case study with qiskit*, IEEE Transactions on Quantum Engineering **3** (2022), pp. 1–17.

[25] X. VASQUES, H. PAIK AND L. CIF, *Application of quantum machine learning using quantum kernel algorithms on multiclass neuron m-type classification*, Scientific Reports **13** (2023), p. 11541.

[26] C.-H. H. YANG ET AL., *A quantum kernel learning approach to acoustic modeling for spoken command recognition*, in ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2023, pp. 1–5.

[27] S. L. WU ET AL., *Application of quantum machine learning using the quantum kernel algorithm on high energy physics analysis at the lhc*, Physical Review Research **3** (2021), p. 033221.

[28] J.-S. KIM, A. MCCASKEY, B. HEIM, M. MODANI, S. STANWYCK AND T. COSTA, *Cuda quantum: The platform for integrated quantum-classical computing*, in 2023 60th ACM/IEEE Design Automation Conference (DAC), IEEE, 2023, pp. 1–4.

[29] R. WILLE, R. VAN METER AND Y. NAVEH, *Ibm's qiskit tool chain: Working with and developing for real quantum computers*, in 2019 Design, Automation & Test in Europe Conference & Exhibition (DATE), IEEE, 2019, pp. 1234–1240.

[30] V. BERGHOLM ET AL., *Pennylane: Automatic differentiation of hybrid quantum-classical computations*, arXiv preprint arXiv:1811.04968 (2018).

[31] S. V. ISAKOV ET AL., *Simulations of quantum circuits with approximate noise using qsim and cirq*, arXiv preprint arXiv:2111.02396 (2021).

[32] J. MARTYN, G. VIDAL, C. ROBERTS AND S. LEICHENAUER, *Entanglement and tensor networks for supervised image classification*, arXiv preprint arXiv:2007.06082 (2020).

[33] T. SUZUKI, T. MIYAZAKI, T. INARITAI AND T. OTSUKA, *Quantum ai simulator using a hybrid cpu–fpga approach*, Scientific Reports **13** (2023), p. 7735.

[34] R. SHAYDULIN ET AL., *Evidence of scaling advantage for the quantum approximate optimization algorithm on a classically intractable problem*, arXiv preprint arXiv:2308.02342 (2023).