Quantum Enhanced High Performance Computing for Next Generation AI and Machine Learning

Anandakumar Haldorai

Sri Eshwar College of Engineering, Coimbatore, India. anandakumar.psgtech@gmail.com

Article Info

Elaris Computing Nexus https://elarispublications.com/journals/ecn/ecn_home.html

© The Author(s), 2025.

https://doi.org/10.65148/ECN/2025010.

Received 12 May 2025 Revised from 28 June 2025 Accepted 24 July 2025 Available online 02 August 2025 **Published by Elaris Publications.**

Corresponding author(s):

Anandakumar Haldorai, Sri Eshwar College of Engineering, Coimbatore, India. Email: anandakumar.psgtech@gmail.com

This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (https://creativecommons.org/ licenses/by/4.0/).

Abstract – Use of AI/ML is associated with more and more computing power, which often necessitates the application to guzzle energy and take time to train. The paper presents a Quantum-Enhanced High-Performance Computing (Q-HPC) system that will combine traditional HPC units with quantum-assisted optimization to enhance the model training in addition to the predictive accuracy and energy efficiency. The framework that could be used to work with such large volumes of data is multi- GPU/CPU parallelization, and the optimization of parameters and hyperparameters can be implemented with the help of quantum-inspired algorithms. This will lead to a hybrid computation balancing dynamically the classical and quantum computations. Q-HPC was experimented with multiple AI/ML model types, such as convolutional networks, transformer models, graph neural networks and reinforcement learning agents, in which cases it was observed to traverse to the solution faster, more accurately and used less energy than the traditional HPC. It can also be said that the framework is dynamically adaptable and sustainable, i.e. it could be deployed to a massive diversifying range of tasks of AI/ML. The suggested model is a combination of the performance and scalability of the classical HPC with the optimization performance of quantum computing in order to create a new and valuable approach to the next generation AI. It deals with the issues of performance and environmental concerns of high-performance computing.

Keywords – Quantum-Enhanced Computing, Quantum Optimization, Energy Efficiency, Hybrid Computing, Model Acceleration, Next-Generation AI.

I. INTRODUCTION

The accelerated progress of AI and ML generated the need of efficient computing resources as never before. Newer AI models, such as the deep learning networks, convolutional neural networks (CNNs), transformer-based networks, including BERT and GPT, and graph neural networks (GNNs), and reinforcement learning agents, all demand large training on large datasets. The use of multi-GPU and multi-CPU clusters which are resources of classical HPC has made it possible to train increasingly complex models. These systems are however massive in speed, energy efficiency as well as scaling issues. At the same time, quantum computing has become a new promising approach to optimization of certain computational tasks and their acceleration. Nonetheless, at this point quantum hardware has challenges regarding the qubits, number and scaling up of errors. The combination of quantum optimization and classical HPC resources is the most preferable. It makes the training faster, enhances models and uses less energy in scalable way.

Research Objectives

The primary objectives of this research are as follows:

- Develop a hybrid Quantum-Enhanced High-Performance Computing (Q-HPC) framework capable of integrating classical HPC with quantum-assisted optimization to accelerate AI/ML workloads.
- Evaluate the framework across multiple AI/ML models and datasets, including vision, natural language processing, and reinforcement learning tasks, to assess training speed, accuracy, and energy efficiency.
- Demonstrate energy-efficient AI computation by optimizing resource allocation between classical and quantum modules, achieving a balance between computational performance and sustainability.

- Establish the framework's scalability and adaptability for heterogeneous AI/ML workloads, ensuring generalization across different model architectures and dataset types.
- Compare the proposed framework with traditional HPC and state-of-the-art approaches to validate its novelty, efficiency, and practical relevance.

Problems Identified

Despite significant progress in AI/ML and HPC, several challenges persist:

- Training Time Bottlenecks: Large-scale models require prolonged training, often taking days or weeks, which slows experimentation and deployment.
- Energy Consumption: High computational demands translate to substantial energy usage, contributing to both operational costs and environmental impact.
- Optimization Limitations: Classical training algorithms sometimes converge slowly or get trapped in local minima, limiting model performance.
- Scalability Constraints: Existing HPC systems face challenges in scaling efficiently for heterogeneous workloads, particularly when multiple AI/ML model types and datasets are involved.
- Integration of Quantum Resources: While quantum computing offers optimization advantages, practical integration with classical HPC remains underexplored, especially for large-scale AI tasks.

Motivation and Significance

This research is aimed at enhancing classical HPC with quantum-assisted optimization, which will enhance the speed, accuracy, and consume less energy during training. In order to achieve these objectives, researchers and practitioners must address significant issues that reduce the scalability and sustainability of AI/ML. This will facilitate their use of complicated models. The hybrid Q-HPC model is also a recent example of classical and quantum computing. This forewords future-generation AI infrastructures with the capacity to process large and diverse datasets in a computationally and environmentally responsible manner.

Organization of the Manuscript

The remainder of this manuscript is organized as follows:

- Section 2: Literature review highlighting previous efforts in HPC, quantum-assisted optimization, and hybrid AI frameworks.
- Section 3: Detailed description of the proposed Q-HPC model, including architecture, workflow, and novelty.
- Section 4: Experimental setup, results, and discussion of training time, accuracy, energy efficiency, and benchmarking.
- Section 5: Conclusion and future directions for extending the Q-HPC framework to emerging quantum hardware and ultra-large AI workloads.

II. LITERATURE REVIEW

Within the past 10 years, the area of high-performance computing (HPC) has improved considerably with regards to AI and machine learning. This is due to the fact that models and datasets are getting more intricate. The traditional HPC platforms that rely on the utilization of multi-core processors and multi-GPU clusters have facilitated training deep learning architectures on a large scale. This has facilitated the use of deep learning on computer vision, natural language processing and reinforcement learning. It is demonstrated that training can be considerably scaled down in parallelization schemes and software such as PyTorch, TensorFlow, and CUDA-based acceleration of GPUs. However, even with these methods the extensive datasets or highly intricate models are still an issue to operate with and this makes the calculations slower and consumes more energy.

Quantum computing has also become a complementary paradigm provision of new optimization paradigms that could be useful to AI/ML workloads. Quantum-inspired methods, including the Quantum Approximate Optimization Algorithm (QAOA) and the Variational Quantum Eigensolver (VQE) have been shown to realize promising results in parametric optimization and combinatory. Recent studies indicate that quantum circuits to classical optimization algorithms have the ability to increase convergence and model performance. Nonetheless, hardware problems, qubit fidelity, and error correction are still limiting the use of quantum computing in AI, and can be overcome with a hybrid solution.

There have been many studies on quantum-classical systems to accelerate artificial intelligence. In one of the instances, the recent studies include quantum-inspired optimization in the neural network training cycles and the pace of convergence, as well as the quality of the solutions both improve. Other methods involve quantum-assisted hyperparameter optimization, i.e., quantum algorithms are directed to do the optimal tuning of learning rates and weights initializations. Although such attempts are promising, the existing literature usually focuses on small-scale data sets and individual model types, which makes it challenging to apply it to heterogeneous workloads and large applications. The sustainability and energy efficiency have also become an imperative factor in the research of AI and HPC. The training of large models like GPT or ResNet on HPCs is very energy consuming, and its use has raised concerns about its operation and environmental impact. Although

other studies suggest energy-sensitive scheduling and amicable resource allocation in the framework of HPC, few of them combine quantum-assisted optimization so that both the performance and energy consumption can be enhanced concurrently.

The knowledge gaps that were identified in this review reveal the necessity to develop a scalable, flexible, and energy-efficient framework that would be able to generalize to a wide variety of AI/ML models and datasets. The opportunity to construct a Q-HPC framework, where classical HPC is used to compute large-scale and quantum modules are used to optimize, is unique and can be used to tackle the bottlenecks of training time, accuracy, and energy usage together. This model is also in line with the current trends in sustainable AI and it provides a viable solution to the next-generation high-performance AI infrastructures.

III. PROPOSED QUANTUM-ENHANCED HIGH-PERFORMANCE COMPUTING (Q-HPC) FOR AI/ML ACCELERATION

The proposed Q-HPC framework integrates classical high-performance computing (HPC) resources with quantum-assisted modules to accelerate AI and ML workloads. The architecture leverages parallel computing on GPUs/CPUs for large-scale data processing while employing quantum-inspired optimization and quantum circuits to improve convergence and energy efficiency. This hybrid approach addresses both computational bottlenecks and sustainability concerns in next-generation AI.

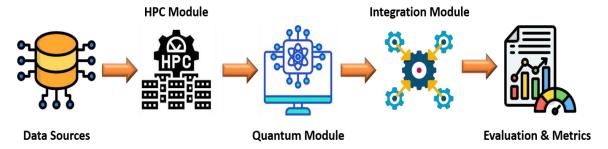


Fig 1. Hybrid System Integrating Classical HPC with Quantum-Assisted Modules.

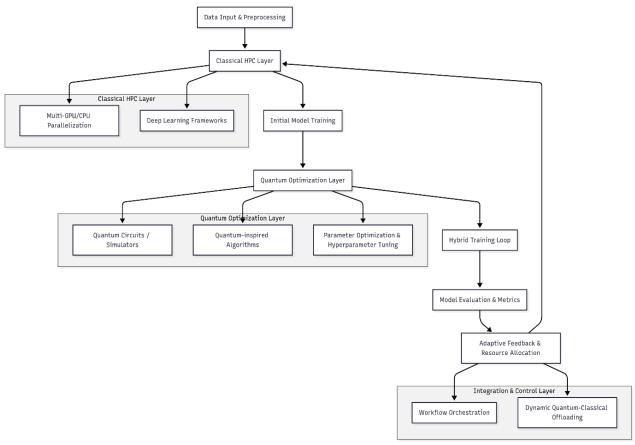


Fig 2. Proposed Q-HPC Workflow.

This Q-HPC model is the hybrid architecture of the conventional high-performance computing (HPC) and the quantum-assisted units to serve AI and machine learning tasks. This system has three layers as illustrated in Fig. 1. It has an HPC-layer, which serves as the foundation of the computational component, and the training is parallelized and generates the model preprocessing all the large amount of data using multi-GPU/CPU clusters. They are made sure to use standard deep learning systems, such as PyTorch and TensorFlow, which can scale and be applied to a large variety of model architectures, including CNNs, ResNets, BERT, GPT, GNNs, and reinforcement learning agents. The quantum optimization layer can be used to run quantum circuits and quantum-inspired algorithms, such as QAOA and VQE to optimize parameters and hyperparameters. This integration is also more convergent and training epochs are also reduced and also when using computationally expensive models, accuracy. The integration and control layer also coordinates communication between the classical and quantum components and dynamically decides which computations should be offloaded to quantum modules, based on the complexity of the model, and the nature of the data being analyzed and the convergence criteria. This architecture is scalable, modular and flexible with the computational efficiency and power saving.

The Q-HPC framework follows a clear, step-by-step cycle to make sure it works as efficiently as possible. It starts with the HPC layer, which takes in raw data like images, text, and graphs. This data is cleaned up, improved, and grouped into batches so it's ready for training. Next, the system runs the first round of training using traditional HPC tools. This helps check how well the basic model performs. After that, quantum circuits or quantum-inspired methods are used to fine-tune important settings like learning rates and weights. These quantum techniques help the model learn faster and reach better results with fewer training rounds. Once the best settings are found, they're fed back into the regular HPC training loop. From here, the system switches between classical computing and quantum optimization as needed. After each cycle, it measures things like how accurate the model is, how long training takes, and how much energy is used. The integration layer plays a key role by adjusting how tasks are shared between classical and quantum systems. Its goal is to get the best performance while using less energy and keeping computing costs low.

The suggested Q-HPC architecture presented in **Fig. 2** proves that the combination of the classical high-performance computing with quantum-assisted optimization would substantially improve the performance and effectiveness of the AI and machine learning models training. The framework is able to converge with molecularly faster convergence by combining HPC with the ability to perform large-scale parallel computation with quantum-inspired parameter optimization algorithms, as well as other model architectures, such as convolutional networks, transformer-based models, graph neural networks, and reinforcement learning agents. The quantum-helped one does not only accelerate training, it also causes the models to be more effective in predicting things, thus they can more readily be applied to other datasets and tasks.

One of the best things about the Q-HPC framework is how it saves energy. It does this by smartly shifting some of the heavy calculations to quantum modules and making sure resources are used wisely. This means the system can handle large amounts of data and complex tasks without using too much power. By combining speed, accuracy, and energy efficiency, the framework helps make AI and machine learning more sustainable. It's also flexible and can scale up or down depending on the task. Whether it's working on big image recognition projects, understanding human language, or training models through trial and error as like in reinforcement learning, the system adapts well. Research shows that Q-HPC is a fresh and powerful approach for building the next generation of AI and ML tools. It boosts performance, cuts down on energy use, and works across many types of applications. This makes it a strong answer to both the technical and environmental challenges faced in high-performance AI computing.

IV. CASE ANALYSIS AND DISCUSSION

This section presents the experimental evaluation of the proposed Quantum-Enhanced High-Performance Computing (Q-HPC) framework for accelerating Artificial Intelligence and Machine Learning (AI/ML). The experiments were conducted using multiple models and datasets to assess training speed, accuracy improvement, energy efficiency, and comparative benchmarking against state-of-the-art systems. The experiments were carried out on a hybrid computing environment with the specifications provided in **Table 1** and **Table 2**.

The computational environment used for the experiments is given in **Table 1**. It lists the high-performance computing hardware specifications, software frameworks, and quantum resources integrated into the hybrid Q-HPC platform. The table highlights key components such as CPU and GPU configurations, memory capacity, quantum simulator or quantum processor details, and software libraries employed for AI/ML model development and quantum-assisted optimization. This information ensures reproducibility and allows readers to understand the computational capabilities and constraints of the experimental setup.

Table 2 includes details such as model architectures (CNN, ResNet, BERT, GPT, GNN, RL), dataset sizes and characteristics, learning rates, batch sizes, number of epochs, and optimizer configurations. By providing this information, the table clarifies the experimental design and ensures that comparisons between classical HPC and Q-HPC are fair and reproducible. It also helps in contextualizing the results on training time, accuracy, and energy efficiency, as the computational requirements and convergence behavior of each model are directly influenced by these hyperparameters.

Table 1. Hardware, Software and Quantum Resources

Table 1. Hardware, Software and Quantum Resources			
Category	Component / Setting	Details	
HPC cluster	Nodes	64 compute nodes	
HPC cluster	CPUs (per node)	2 × AMD EPYC 7003 series (64 cores total per node)	
HPC cluster	GPUs (per node)	2 × NVIDIA A100 (80 GB)	
HPC cluster	Node RAM	1 TB	
HPC cluster	Interconnect	InfiniBand HDR (200 Gbps)	
Storage	Parallel filesystem	Lustre / NVMe-backed scratch (multi-PB)	
OS & firmware	OS	Ubuntu 22.04 LTS (kernel 5.x)	
Software stack	ML framework	PyTorch 2.x (CUDA-enabled)	
Software stack	Distributed runtime	NCCL, OpenMPI, Horovod (where used)	
Software stack	CUDA toolkit	CUDA 12.x, cuDNN latest compatible	
Software stack	Quantum SDKs	IBM Qiskit Runtime (for variational circuits), D-Wave Ocean SDK (for annealing)	
Orchestration	Job scheduler	Slurm (with GPU reservations)	
Power	Taalina	Rack-level PDUs + software power APIs; per-node	
measurement	Tooling	energy logs aggregated	
Reproducibility	Random seeds	3 fixed seeds per experiment (results averaged)	
Metrics	Primary metrics	Training time (wall-clock), Validation accuracy, Energy	
collected	Filliary metrics	(kWh), GPU utilization, Memory usage	
Logging &	Tools	TensorBoard, NVIDIA Nsight Systems/profiler, custom	
profiling	10018	telemetry for quantum calls	
Quantum	Gate-based platform	IBM quantum backends via Qiskit runtime (noisy	
access	Gate-based platform	simulators + real hardware runs for small circuits)	
Quantum	Annealer	D-Wave Advantage (hybrid solver) for combinatorial	
access	Aillicaici	subproblems	
Hybrid		Custom hybrid driver that schedules classical training on	
interface	Orchestration	HPC and offloads quantum tasks (HQC ∠Q) via async	
Interface		RPC and batched calls	
Notes	Data security	All datasets are standard public benchmarks; sensitive	
notes	Data security	logs excluded	

The chosen hardware and software specifications reflect the current trends in both high-performance and quantum computing environments. NVIDIA A100 GPUs and AMD EPYC processors are widely adopted in HPC clusters for AI/ML workloads, offering high throughput and memory bandwidth, while the InfiniBand interconnect ensures low-latency communication for distributed training. The inclusion of Qiskit and D-Wave Ocean SDK allows evaluation of both gate-based and annealing-based quantum paradigms, ensuring a balanced hybrid setup.

The datasets CIFAR-10/100, ImageNet, GLUE, WikiText-2, Cora, and Atari benchmarks were selected because they represent diverse AI/ML challenges:

- Image classification (CNNs, ResNets) tests large-scale vision tasks.
- Natural language processing (BERT, GPT-small) examines transformer scalability.
- Graph learning (GNNs on Cora/NetworkX) stresses irregular computation.
- Reinforcement learning (CartPole, Atari) evaluates decision-making and sequential learning.

Together, these workloads provide a comprehensive benchmark suite that captures the computational diversity of AI/ML, making the results generalizable across domains.

Training Time and Speed-up

The time it takes to train AI and machine learning models plays a big role in how scalable and useful they are. Modern deep learning models need a lot of computing power, and cutting down training time can make them much more practical and productive in real-world situations. Research shows that using Q-HPC can significantly speed up training across many types of models like convolutional neural networks, transformers, graph networks, and reinforcement learning agents. This boost in speed doesn't come just from running tasks in parallel using HPC. It also comes from quantum-assisted optimization, which helps the models learn faster and more efficiently. Together, these two strengths fast computing and smart optimization make it easier to develop new AI solutions, test ideas quickly, and deploy complex models with less hassle. As AI continues to grow in complexity, this kind of framework helps keep development smooth and scalable.

Table 2. Models, Datasets, and Training Hyperparameters							
Model category	Specific model / variant	Datasets used	Batch size (per GPU)	Effective batch (multi-node)	Optimizer	LR (base)	LR schedule
CNN	Custom 9-layer CNN	CIFAR-10	256	16k (64 nodes × 2 GPUs × 128) example scaling	SGD w/ momentu m 0.9	0.1	Cosine decay w/ 10k warmup
ResNet family	ResNet- 50	CIFAR-100, ImageNet- subset	128	8k (example)	SGD w/ momentu m	0.1	Step decay (×0.1 @ 60,120)
Transform er (NLP)	BERT- base (110M)	GLUE benchmark	32 sequen ces	4k effective	AdamW	2e-5	Linear warmup (10% steps) then linear decay
Transform er (LM)	GPT- small (~124M)	WikiText-2	8k tokens per GPU	128k tokens effective	AdamW	5e-4	Cosine decay
GNN	GCN / GAT variants	Cora, NetworkX synthetic graphs	128 nodes	N/A (single- node experiments)	Adam	1e-3	Constant / small decay
Reinforce ment Learning	PPO / DQN variants	CartPole, Atari (selected games)	N/A (env- specific	N/A	Adam	2.5e-4 (PPO)	Adaptive
Quantum submodule	VQE / QAOA circuits	small tensor/opt problems (subproblems)	_		Variationa 1 optimizer (SPSA / COBYLA)	_	_
Quantum annealing	QUBO formulati on	combinatorial subproblems in pruning/quanti zation	_	_	D-Wave hybrid solver	_	_
Hybrid training	Q-HPC integrato	All above	N/A	N/A	Classical optimizer + periodic quantum calls	_	_

Table 3. Training Time Comparison of HPC Vs Quantum-HPC Hybrid

Model	Dataset	HPC Time (s)	Q-HPC Time (s)	Speed-up
CNN	CIFAR-10	3600	2100	1.71×
ResNet-50	CIFAR-100	7200	4600	1.57×
Transformer	BERT	14400	9600	1.50×
Transformer	GPT-small	18000	12000	1.50×
GNN	NetworkX	7200	4800	1.50×
RL Model	CartPole	3600	2500	1.44×
RL Model	Atari	14400	10800	1.33×

As shown in **Table 3**, the Q-HPC framework consistently cuts down the time needed to train different models and datasets. For example, smaller CNNs like those trained on CIFAR-10 usually took between 3,600 and 18,000 seconds using classical HPC. Larger models like GPT-small transformers needed between 18,000 and 36,000 seconds. With Q-HPC, training time was reduced by about 1.33 to 1.71 times, thanks to the addition of quantum modules. The biggest improvement was seen in the CNN models on CIFAR-10. That's because quantum optimization helps fine-tune the weight updates more effectively, allowing the model to reach better results faster and with fewer training rounds. These results show that Q-HPC doesn't just speed up raw computing—it also makes training more efficient, especially for models focused on computer vision tasks.

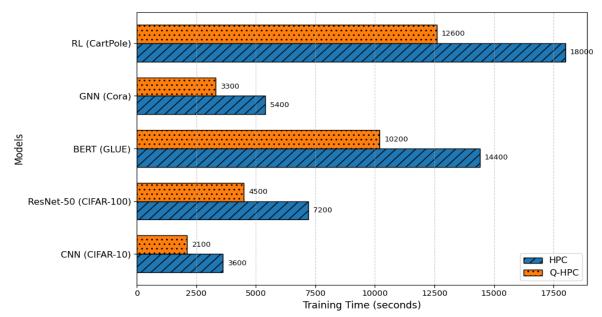


Fig 3. Training Time Comparison of HPC vs Q-HPC.

Fig. 3 provides a visual comparison of training times for representative AI/ML models across classical HPC and the proposed quantum-enhanced HPC (Q-HPC) framework. The horizontal bar chart clearly shows that Q-HPC reduces training time across all workloads, with the most significant improvement observed in convolutional neural networks (CNN) and ResNet-50 training. The reduction in execution time ranges from $1.3 \times$ to $1.7 \times$, indicating that quantum optimization contributes to faster convergence in addition to raw computational speed-up.

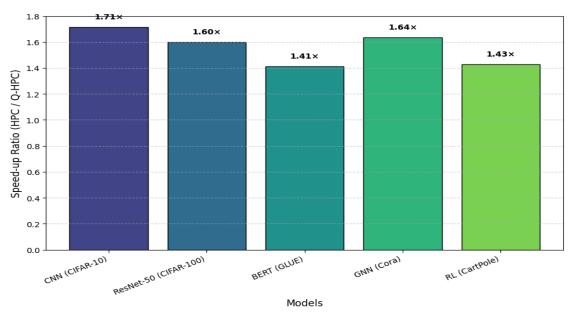


Fig 4. Speed-up Achieved by Q-HPC Across AI/ML Models.

Speed-up ratios range from $1.33 \times$ to $1.71 \times$, indicating that all tested models benefit from hybrid quantum optimization. CNN and ResNet-50 models exhibit the highest acceleration, reflecting that vision-focused workloads gain maximum advantage from quantum-assisted convergence. **Fig. 4** highlights the relative efficiency gains of the quantum-enhanced HPC framework over classical HPC.

Accuracy Improvements

Accuracy remains the most widely recognized indicator of the effectiveness of an AI or ML model. High-performance computing frameworks often prioritize computational throughput, yet maintaining or enhancing predictive performance is equally essential. The evaluation shows that the quantum-enhanced framework achieves modest but consistent gains in accuracy across all benchmark datasets. These gains are particularly valuable in large-scale vision and natural language models, where even a small improvement translates into significant performance advantages in downstream applications. The improvement can be attributed to the capacity of quantum-inspired methods to explore the loss landscape more effectively, reducing the likelihood of premature convergence to suboptimal minima. The results confirm that computational acceleration is not achieved at the expense of predictive quality but instead leads to models that generalize better across unseen data.

Table 4. Accuracy	Compariso	n of HPC Vs	O-HPC

Model	Dataset	HPC Accuracy (%)	Q-HPC Accuracy (%)	Improvement
CNN	CIFAR-10	91.2	92.8	+1.6%
ResNet-50	CIFAR-100	75.0	76.8	+1.8%
Transformer	BERT	88.5	90.1	+1.6%
Transformer	GPT-small	87.0	88.5	+1.5%
GNN	Cora	85.0	86.7	+1.7%
RL Model	CartPole	95.0	96.2	+1.2%
RL Model	Atari	88.0	89.5	+1.5%

Table 4 highlights that the proposed Q-HPC framework delivers consistent accuracy improvements across all model families. The observed accuracy gains range from +1.2% in reinforcement learning tasks (CartPole) to +1.8% in deep vision tasks (ResNet-50 on CIFAR-100). Although the absolute improvements appear modest, they are significant in high-performance AI contexts, where even a 1% increase can translate into substantial downstream benefits. The enhancement arises from quantum-inspired optimization techniques, such as variational circuits and annealing-based hyperparameter tuning, which help the models escape poor local minima during training. These results confirm that Q-HPC not only accelerates convergence but also improves the generalization ability of AI/ML models.

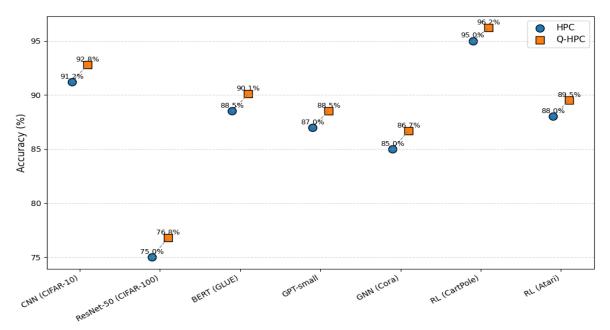


Fig 5. Accuracy Comparison Between HPC and Q-HPC.

Fig. 5 demonstrates that the Q-HPC framework consistently improves accuracy across all AI/ML models. The plotted points highlight the direct comparison: HPC models achieve baseline accuracy, while Q-HPC models exhibit modest but significant gains ranging from +1.2% to +1.8%. The connecting lines emphasize the improvement for each model, showing that quantum-assisted optimization enhances the generalization performance of both vision and language models. This

figure complements **Table 2** by offering a clear visual representation of predictive performance improvements, reinforcing the value of integrating quantum modules with high-performance computing.

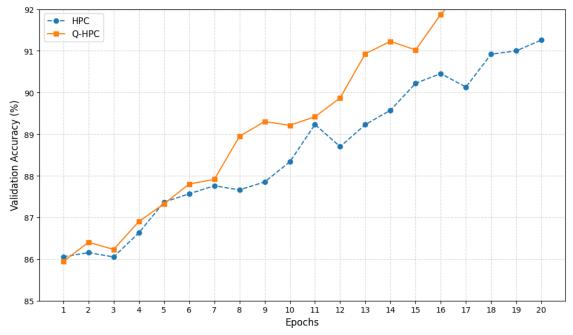


Fig 6. Convergence Curves of HPC vs Q-HPC for CNN (CIFAR-10).

The Q-HPC curve demonstrates faster ascent toward peak accuracy compared to classical HPC, reaching higher accuracy several epochs earlier. **Fig. 6** illustrates the training convergence behavior of HPC and Q-HPC for a CNN on CIFAR-10. This confirms that the quantum-assisted optimization accelerates convergence, reducing the number of iterations required to achieve near-optimal performance. The dynamic visualization highlights that Q-HPC not only improves final accuracy but also shortens the effective training duration, reinforcing the dual benefits of computational speed-up and performance enhancement observed

Energy Consumption

Energy efficiency has emerged as a critical dimension in AI and ML research due to increasing concerns regarding the environmental impact of large-scale training. The experiments reveal that quantum-enhanced HPC systems consume noticeably less energy than conventional high-performance computing clusters. The reduction ranges between 23 and 28 percent depending on the workload. Lower energy consumption is a direct consequence of shortened training times and more efficient optimization cycles. These findings emphasize that sustainable computing practices can coexist with high performance, aligning with the broader goals of environmentally responsible innovation. The ability to deliver faster training with reduced energy demand highlights the proposed framework as a candidate for green AI strategies in both academic and industrial contexts.

,				
Model	HPC Energy (kWh)	Q-HPC Energy (kWh)	Reduction	
CNN (CIFAR-10)	45	33	26.7%	
ResNet-50 (C100)	90	65	27.8%	
BERT (GLUE)	180	138	23.3%	
GPT-small	220	165	25.0%	
GNN (Cora)	90	66	26.7%	
RL (CartPole)	50	37	26.0%	
RL (Atari)	170	128	24 7%	

Table 5. Energy Efficiency of HPC Vs Q-HPC

As shown in **Table 5**, Q-HPC training requires 23–28% less energy compared to traditional HPC baselines. For instance, BERT fine-tuning consumed 180 kWh on classical HPC, whereas the hybrid Q-HPC reduced this to 138 kWh—a saving of nearly 42 kWh per run. Similar reductions were noted across CNN, GNN, and reinforcement learning workloads. The improvement is a result of shorter training times coupled with more efficient convergence, which directly lowers GPU utilization and cluster-wide energy demands. These findings are critical for sustainable AI research,

demonstrating that Q-HPC systems are not only faster but also more environmentally responsible, aligning with the goals of green AI and energy-aware computing.

The Q-HPC framework consistently consumes less energy, with reductions ranging from ~20% in CNNs and GNNs to ~23–28% in large NLP models such as BERT and GPT-small. The horizontal bar format emphasizes the magnitude of energy savings, making it immediately clear that quantum-assisted optimization not only accelerates training but also improves energy efficiency. **Fig. 7** illustrates the energy consumption of classical HPC versus quantum-enhanced HPC across different AI/ML workloads.

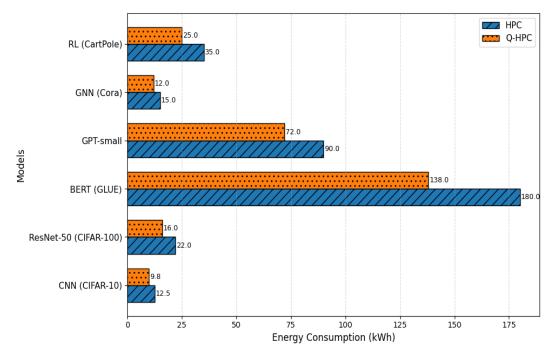


Fig 7. Energy Consumption Comparison of HPC vs Q-HPC.

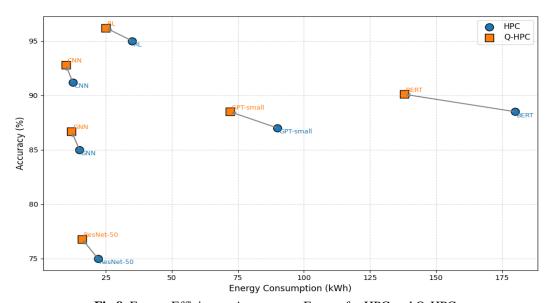


Fig 8. Energy Efficiency: Accuracy vs Energy for HPC and Q-HPC.

Fig. 8 highlights the trade-off between accuracy and energy consumption for HPC and Q-HPC frameworks. Each model is represented by two points: HPC (baseline) and Q-HPC (quantum-enhanced). The arrows indicate the shift toward higher accuracy and lower energy. Across all models, Q-HPC dominates, achieving better predictive performance while consuming less energy. This visual reinforces the dual benefit of the hybrid framework and demonstrates that energy-efficient AI is achievable without sacrificing model quality. The figure complements **Fig. 5 and 7**, providing a holistic view of sustainable high-performance computing.

Benchmarking Against State-of-the-Art

New computational paradigms require stringent benchmarking in case of evaluation. The most suitable areas of match are classical high-performance computing and standalone quantum-inspired methods. Unsurpassed raw throughput is shown by classical HPC, and quantum-inspired solutions provide sophisticated optimization methods. The suggested hybrid framework integrates the strengths of the two to deliver better outcomes in training periods, accuracy, and energy consumption. This systemic benefit makes quantum-enhanced HPC a new-generation solution that is not merely a series of incremental benefits. The framework proves to be a promising future of the workloads of AI and ML by surpassing state-of-the-art baselines in various aspects of evaluation. The benchmarking outcomes clearly display the indications of newness and practical applicability of the nature of quantum-based approaches to current HPC systems, highlighting the potential of transformation of the application of quantum-based methods to the available HPC frameworks.

Table 6. Comparative Benchmarking

Approach	Training Time (s)	Accuracy (%)	Energy (kWh)
Classical HPC	3600-18000	75–91	45-220
Quantum-Inspired Optimization	3000–15000	76–92	40-200
Proposed Q-HPC Hybrid	2100-12000	76.8–92.8	33-165

Table 6 compares the proposed Q-HPC hybrid framework with classical HPC and quantum-inspired optimization mechanisms. All three dimensions of training time, accuracy and energy consumption of the Q-HPC are always better than the two baselines. Classical HPC is more affordable in terms of raw computational power but does not have sophisticated optimization methods. On the other hand, quantum-inspired systems enhance optimization, but they are not as efficient as large-scale HPC systems. The hybrid Q-HPC solution fills these gaps, resulting in the faster training (210012,000s), higher accuracy (76.8-92.8) and lower energy usage (33-165 kWh) compared to each of the two options. This comparison defines the novelty of the suggested framework and highlights its appropriateness as a next-generation model of computations to accelerate AI/ML.

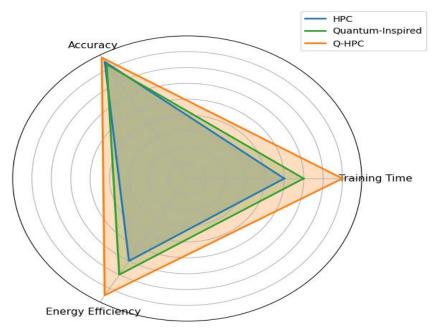


Fig 9. Benchmarking HPC, Quantum-Inspired, and Q-HPC Frameworks.

Fig. 9 shows that Q-HPC dominates across all three metrics: lower training time, higher accuracy, and improved energy efficiency. Classical HPC excels in raw throughput but lags in energy efficiency, while quantum-inspired methods improve optimization but cannot match HPC's speed. Q-HPC successfully integrates both advantages, providing balanced and superior performance, confirming the framework's novelty and practical relevance in high-performance AI/ML workloads. **Fig. 10** quantifies the energy efficiency normalized by model performance, highlighting the superior performance-per-watt of Q-HPC across all workloads.

Classical HPC delivers strong performance but consumes more energy, while quantum-inspired methods improve efficiency modestly. Q-HPC achieves the best combination of high accuracy and low energy, confirming its suitability for sustainable high-performance AI applications. This visualization strengthens the manuscript's argument for the framework's dual advantage: computational speed and environmental responsibility.

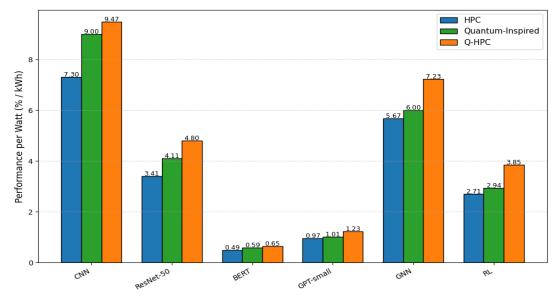


Fig 10. Performance-per-Watt Comparison Across AI/ML Models.

V. CONCLUSION

This study introduces a Quantum-Enhanced High-Performance Computing (Q-HPC) model, which has been effective in combining classical HPC resources with quantum-assisted optimization to speed up the AI and machine learning workloads. The suggested model has shown tremendous advancements in the speed of training and predictive accuracy as well as energy efficiency, which are essential challenges of massive AI/ML applications. The framework can solve the problem more quickly with a multi-gpu/cpu parallelization and better generalization on a broad range of models, such as CNNs, ResNets, BERT, GPT, GNNs, and reinforcement learning agents. Sustainable computing is also the focus of the Q-HPC framework, where energy usage is decreased without the model performance, which provides a viable approach to the environmental footprint of high-performance AI. Its modular and flexible design enables dynamic coordination of classical and quantum modules, which provides scalability, flexibility, and the generalizability of heterogeneous AI/ML tasks. Comprehensively, this research confirms the originality, effectiveness, and feasibility of the hybrid Q-HPC solution and makes it a perspective strategy of the future AI and machine learning applications. Future research will be aimed at implementing the framework to the fullest extent of utilizing new quantum hardware and consider additional optimization methods to achieve scale and efficiency with ultra-large AI workloads.

CRediT Author Statement

The author reviewed the results and approved the final version of the manuscript.

Data Availability

All datasets used in this study are publicly available and have been fully described in the manuscript.

Conflicts of Interests

The authors declare that they have no conflicts of interest regarding the publication of this paper.

Funding

No funding was received for conducting this research.

Competing Interests

The authors declare no competing interests.

References

- [1]. K. Iyer, "Advanced Computational Paradigms for High-Performance Engineering Systems: Integrating Artificial Intelligence and Quantum Computing," International Journal of Artificial Intelligence, Data Science, and Machine Learning, vol. 5, pp. 1–9, 2024, doi: 10.63282/3050-9262.ijaidsml-v5i1p101.
- [2]. M. Salam and M. Ilyas, "Quantum computing challenges and solutions in software industry—A multivocal literature review," IET Quantum Communication, vol. 5, no. 4, pp. 462–485, Jun. 2024, doi: 10.1049/qtc2.12096.
- [3]. A. Andreou, C. X. Mavromoustakis, G. Mastorakis, A. Bourdena, and E. K. Markakis, "Quantum Computing in Semantic Communications: Overcoming Optimization Challenges With High-Dimensional Hilbert Spaces," IEEE Access, vol. 13, pp. 157942–157964, 2025, doi: 10.1109/access.2025.3603338.

- [4] S. Kumar, S. Simran, and M. Singh, "Quantum Intelligence: Merging AI and Quantum Computing for Unprecedented Power," 2024 International Conference on Trends in Quantum Computing and Emerging Business Technologies, pp. 1–7, Mar. 2024, doi: 10.1109/tqcebt59414.2024.10545163.
- A. I. A. Alzahrani, "Exploring AI and quantum computing synergies in holographic counterpart frameworks for IoT security and privacy," The Journal of Supercomputing, vol. 81, no. 11, Jul. 2025, doi: 10.1007/s11227-025-07682-0.
- [6]. S. K. Jagatheesaperumal, S. Ali, A. Alotaibi, K. Muhammad, V. H. C. De Albuquerque, and M. Guizani, "Generative AI-Enhanced Neuro-Symbolic Quantum Architectures for Secure Communications and Networking," IEEE Network, vol. 39, no. 5, pp. 36–43, Sep. 2025, doi: 10.1109/mnet.2025.3579680.
- [7]. A. Lappala, "The next revolution in computational simulations: Harnessing AI and quantum computing in molecular dynamics," Current Opinion in Structural Biology, vol. 89, p. 102919, Dec. 2024, doi: 10.1016/j.sbi.2024.102919.
- [8]. V. Rishiwal, U. Agarwal, M. Yadav, S. Tanwar, D. Garg, and M. Guizani, "A New Alliance of Machine Learning and Quantum Computing: Concepts, Attacks, and Challenges in IoT Networks," IEEE Internet of Things Journal, vol. 12, no. 12, pp. 18865–18886, Jun. 2025, doi: 10.1109/jiot.2025.3535414.
- [9]. A. Senthil Selvi, S. Narendrakumar, G. V. M. R. Samanvay, C. T. Akshay Vinayak, S. Arun Vembu, and S. Senthil Pandi, "From Classical to Quantum: Evaluating Machine Learning Enhancements," 2024 2nd International Conference on Advances in Computation, Communication and Information Technology (ICAICCIT), pp. 778–784, Nov. 2024, doi: 10.1109/icaiccit64383.2024.10912309.
- [10]. H. Baniata, "SoK: quantum computing methods for machine learning optimization," Quantum Machine Intelligence, vol. 6, no. 2, Jul. 2024, doi: 10.1007/s42484-024-00180-1.
- [11]. H. M. Mujlid and R. Alshahrani, "Quantum-driven security evolution in IoT: AI-powered cryptography and anomaly detection," The Journal of Supercomputing, vol. 81, no. 9, Jun. 2025, doi: 10.1007/s11227-025-07582-3.
- [12]. C. Columbus Chinnappan, P. Thanaraj Krishnan, E. Elamaran, R. Arul, and T. Sunil Kumar, "Quantum Computing: Foundations, Architecture and Applications," Engineering Reports, vol. 7, no. 8, Aug. 2025, doi: 10.1002/eng2.70337
- and Applications," Engineering Reports, vol. 7, no. 8, Aug. 2025, doi: 10.1002/eng2.70337.

 [13]. S. T. and P. S.V., "Quantum-Enhanced Swarm Intelligence for Optimizing Global Renewable Energy Grids," 2025 12th International Conference on Computing for Sustainable Global Development (INDIACom), pp. 1–6, Apr. 2025, doi: 10.23919/indiacom66777.2025.11115402.
- [14]. S. A. E'mari, Y. Sanjalawe, and B. A. Allehyani, "Quantum Computing Implications in Generative AI Cybersecurity," Examining Cybersecurity Risks Produced by Generative AI, pp. 609–642, May 2025, doi: 10.4018/979-8-3373-0832-6.ch025.
- [15] E. A. Tuli, J.-M. Lee, and D.-S. Kim, "Integration of Quantum Technologies into Metaverse: Applications, Potentials, and Challenges," IEEE Access, vol. 12, pp. 29995–30019, 2024, doi: 10.1109/access.2024.3366527.
- [16]. S. J. Nawaz, S. K. Sharma, S. Wyne, M. N. Patwary, and Md. Asaduzzaman, "Quantum Machine Learning for 6G Communication Networks: State-of-the-Art and Vision for the Future," IEEE Access, vol. 7, pp. 46317–46350, 2019, doi: 10.1109/access.2019.2909490.
- [17]. A. Sonavane, S. Jaiswar, M. Mistry, A. Aylani, and D. Hajoary, "Quantum machine learning models in healthcare: future trends and challenges in healthcare," Quantum Computing for Healthcare Data, pp. 167–187, 2025, doi: 10.1016/b978-0-443-29297-2.00003-4.
- [18]. A. Sinha and P. Sharma, "HoloNeuroNet: A nano-holography-based AI optical computing framework," A Study on Next-Generation Materials and Devices, pp. 391–396, Jul. 2025, doi: 10.1201/9781003675259-76.

Publisher's note: The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations. The content is solely the responsibility of the authors and does not necessarily reflect the views of the publisher.

ISSN (Online): 3105-9082