Defining the Accelerated Quantum Supercomputer

Naruhiko Tan, NVIDIA



Quantum Error Correction















Year of Quantum Breakthroughs











New Architectures













The Accelerated Quantum Supercomputer

• Al supercomputing architecture connecting quantum hardware

- Hybrid applications combining HPC, Al and QC resources
- A software platform enabling domain scientists
- Qubit-agnostic development of control and error correction







CUDA-Q **Open-Source Software Platform for Accelerated Quantum** Supercomputers

The NVIDIA Quantum Platform Bringing AI Supercomputing to Enable Useful Quantum Computing









- Anyon Technologies
- () Infleqtion
- |QM|
- OQC





















CUDA-Q **Open-Source Software Platform for Accelerated Quantum** Supercomputers

Respectively and the second se New at GTC25

QUANTINUUM Free Access to H1-1 Through CUDA-Q

The NVIDIA Quantum Platform Bringing AI Supercomputing to Enable Useful Quantum Computing







- Anyon Technologies
- (C) Inflection
- |QM|
- OQC

















The NVIDIA Quantum Platform Bringing AI Supercomputing to Enable Useful Quantum Computing

CUDA-Q **Open-Source Software Platform for Accelerated Quantum** Supercomputers

cuQuantum Accelerated Quantum Simulations



The NVIDIA Quantum Platform Bringing Al Supercomputing to Enable Useful Quantum Computing

CUDA-QX Turnkey Quantum Research and Application Development

CUDA-Q **Open-Source Software Platform for Accelerated Quantum** Supercomputers

cuQuantum Accelerated Quantum Simulations

DGX Quantum Reference Architecture for Low Latency Integration

QUANTUM MACHINES

The NVIDIA Quantum Ecosystem Accelerating the Quantum World

Challenges Ahead How to Turn Qubits into Accelerated Quantum Supercomputers

Developing Better QPUs with Simulation

Google Quantum Al

Noise Limits Today's Quantum Hardware

Current quantum hardware is limited to just hundreds of instructions by noise

IVIDIA,

- 2

Quantum Processor

Qiskit Dynamics (CPU)

CPU: Intel Xeon 8480CL | GPU: H100 Energy levels: 64 Transmon x 128 Resonator x 4 Purcell Filter

Dynamics in CUDA-Q Enabling Better QPU Designs

Transmon Simulation (Qubit, Resonator, Filter) Energy Levels = (64, 256, 4)

Time(s)

44 minutes

Qiskit Dynamics (CPU)

CUDA-Q (DGX H100) 0.3s

CPU: Intel Xeon 8480CL | GPU: H100 Energy levels: 64 Transmon x 128 Resonator x 4 Purcell Filter

Dynamics in CUDA-Q Enabling Better QPU Designs

Transmon Simulation (Qubit, Resonator, Filter) Energy Levels = (64, 256, 4)

Time(s)

44 minutes

Better Algorithms with Large-Scale Simulations and Al

Better Algorithms with Large-Scale Simulations and Al

IonQ/AWS/AstraZeneca **AF-QMC** simulation for electronic structure modeling

https://arxiv.org/html/2411.10406v1 (2024) GTC Talk DD73669

HPE Labs **Distributed QC with** Adaptive Circuit Knitting

NCHC 784 qubit simulation to validate QML approaches

© NVIDIA.

Better Algorithms with Large-Scale Simulations and Al

https://arxiv.org/html/2405.02630v2 (2024)

© NVIDIA.

10^{-3} State of the art error rates $< 10^{-10}$

Expected Error rates needed

10^{-3} State of the art error rates < 10^{-10} Expected Error rates needed

Gidney, Ekera. How to factor 2048 bit RSA integers in 8 hours using 20 million noisy qubits. https://arxiv.org/abs/1905.09749 (2019)

14,238

Qubits needed to factor a 2,048 bit number in 8 hours

22,325,184

Qubits needed with QEC

100TB/s

Data streaming from 1M qubits

<mark> NVIDIA</mark>.

Decoding Time

O(Gate Time)

Code Distance

Practical QEC

Decoding Time

O(Gate Time)

Code Distance

Practical QEC

Decoding Time

O(Gate Time)

Code Distance

Parallelization

Practical QEC

Quantum Error Correction – Why Quantum Computing Needs AI Supercomputing

Decoding Time

O(Gate Time)

Code Distance

Practical QEC

<mark> NVIDIA</mark>.

 Trained on synthetic data with QuEra noise model

 Improved accuracy over MLE for small codes - d = 3

 Decoding time ~1ms, 50x improvement over SotA

Decoding QuEra's QPU with Al Transformer-based decoder

Future Work

 Larger codes – much more training data needed

 Fine-tune on QPU data for higher accuracy

Deploy for real-time decoding

Decoding QuEra's QPU with Al Transformer-based decoder

 $\alpha = 3$

$$d=11$$

Code Distance

Custom Noise Model

Sampling Algorithm

Building an Al Decoder

Strategic Kraus **Operator Sets**

> $\{K_0, ..., K_N\}$ $\{K_0, ..., K_N\}$ {K₀,..., K_N}

Synthetic Training Data

Custom Noise Model

Sampling Algorithm

On DGX H100 Supercomputer with 576 nodes:

Released open source in CUDA-Q 0.10

Building an Al Decoder

Strategic Kraus **Operator Sets**

Emulated 85 qubit system – 1 million shots generated in > 6 hours Emulated 35 qubit system – 1 trillion shots generated in > 2 hours

Accelerated Simulators

Synthetic Training Data

<mark> NVIDIA</mark>

Announcing DGX Quantum Alpha Release

Integrated system for low-latency GPU-QPU programming

Building block for accelerated quantum supercomputer

<a>4us latency and real-time RL control demonstrated

Now shipping in Alpha

QUANTUM MACHINES

Supercomputers are Driving Quantum Computing Progress

AIST - Japan QuEra, Fujitsu, OptQC

Novo Nordisk Foundation - Denmark

PSNC - Poland ORCA Computing

Julich - Germany IQM, D-Wave

GENCI - France Pasqal, Quandela

NQCC – UK **ORCA Computing, SEEQC**

Pawsey - Australia **Quantum Brilliance**

LRZ – Germany AQT, IQM, planqc

© NVIDIA.

Thank You

