

# Quantum Accelerated Supercomputing

A Step-by-Step Guide to Expanding Simulation  
Capabilities and Enabling Interoperability of Quantum Hardware

Monica VanDieren  
NVIDIA

Jin-Sung Kim  
NVIDIA

Ermal Rrapaj  
NERSC at LBNL

Katie Klymko  
NERSC at LBNL

Zohim Chandani  
NVIDIA

Tommaso Macri  
QuEra

# Agenda

## Morning Sessions

- Introductions
- Overview of CUDA-Q
- Onboarding to Perlmutter
- Accelerating Quantum Simulations Notebook
- NVIDIA Quantum Cloud and Local install of CUDA-Q

## Afternoon Sessions

- Intro to large scale clusters, how to navigate and use them
- Example: Quantum chemistry at NERSC
- Industry Use Case: Simulating Hamiltonians of molecules with 30,000 terms
- QuEra's Quantum Hardware
- Example: Quantum reservoir computing

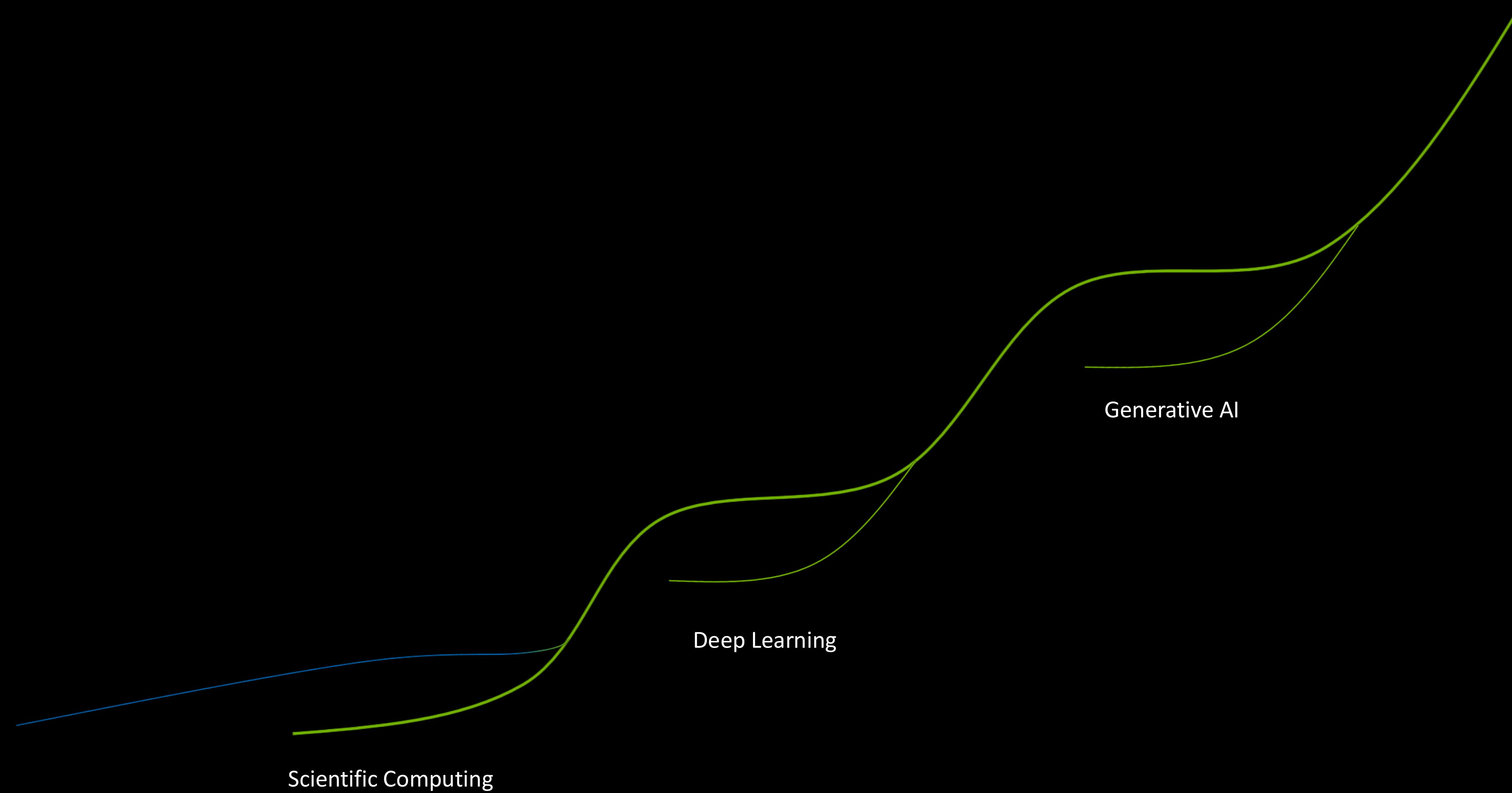
# CUDA-Q Overview



# Agenda

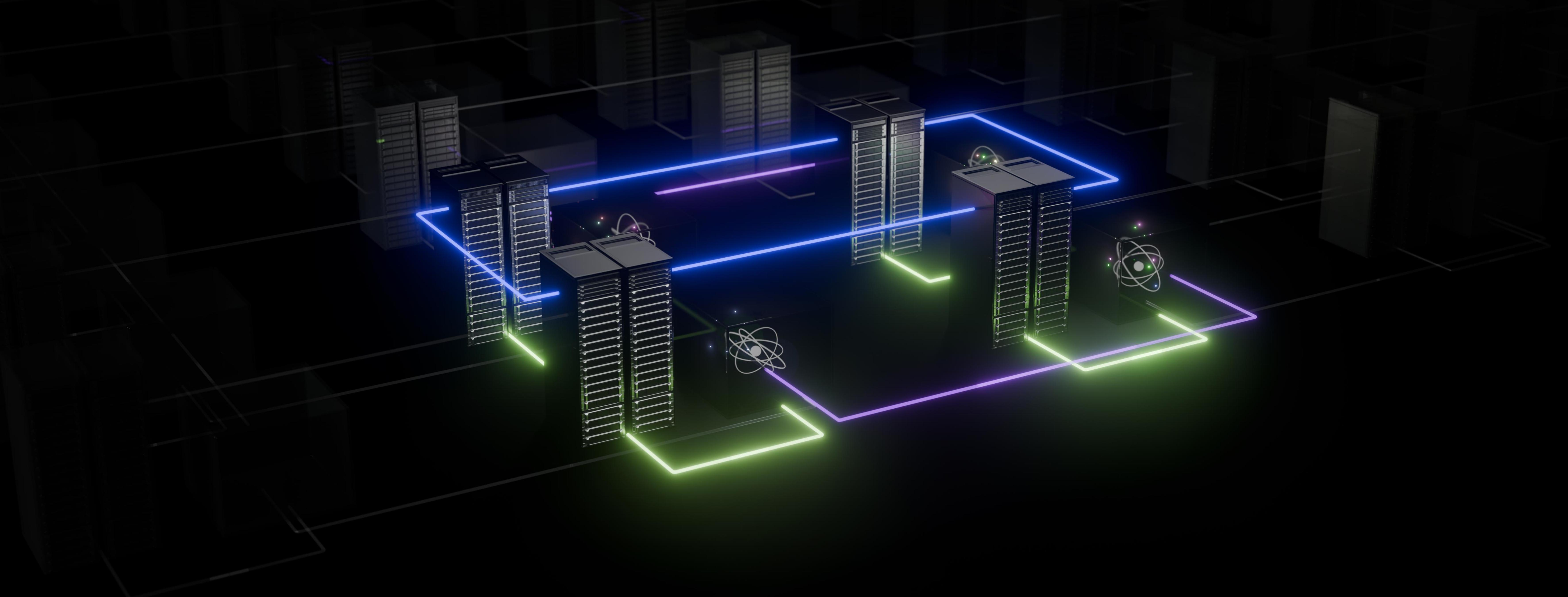
- Accelerated Quantum Supercomputer
- CUDA-Q Platform
- Accelerating Quantum Simulations with GPUs
- Installing CUDA-Q and the NVIDIA Quantum Cloud
- CUDA-Q Resources

# NVIDIA's History of Enabling Computing Revolutions



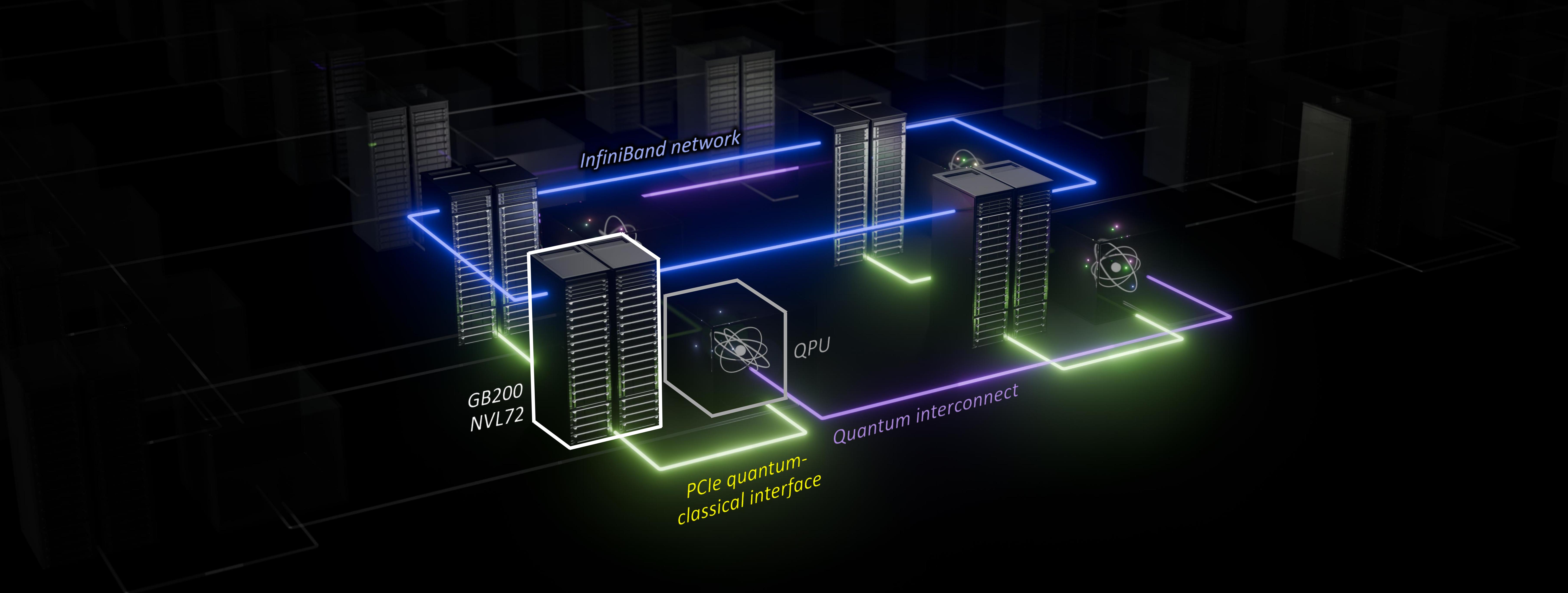
# NVIDIA's History of Enabling Computing Revolutions





# DGX SuperPOD integrates the first useful QPUs

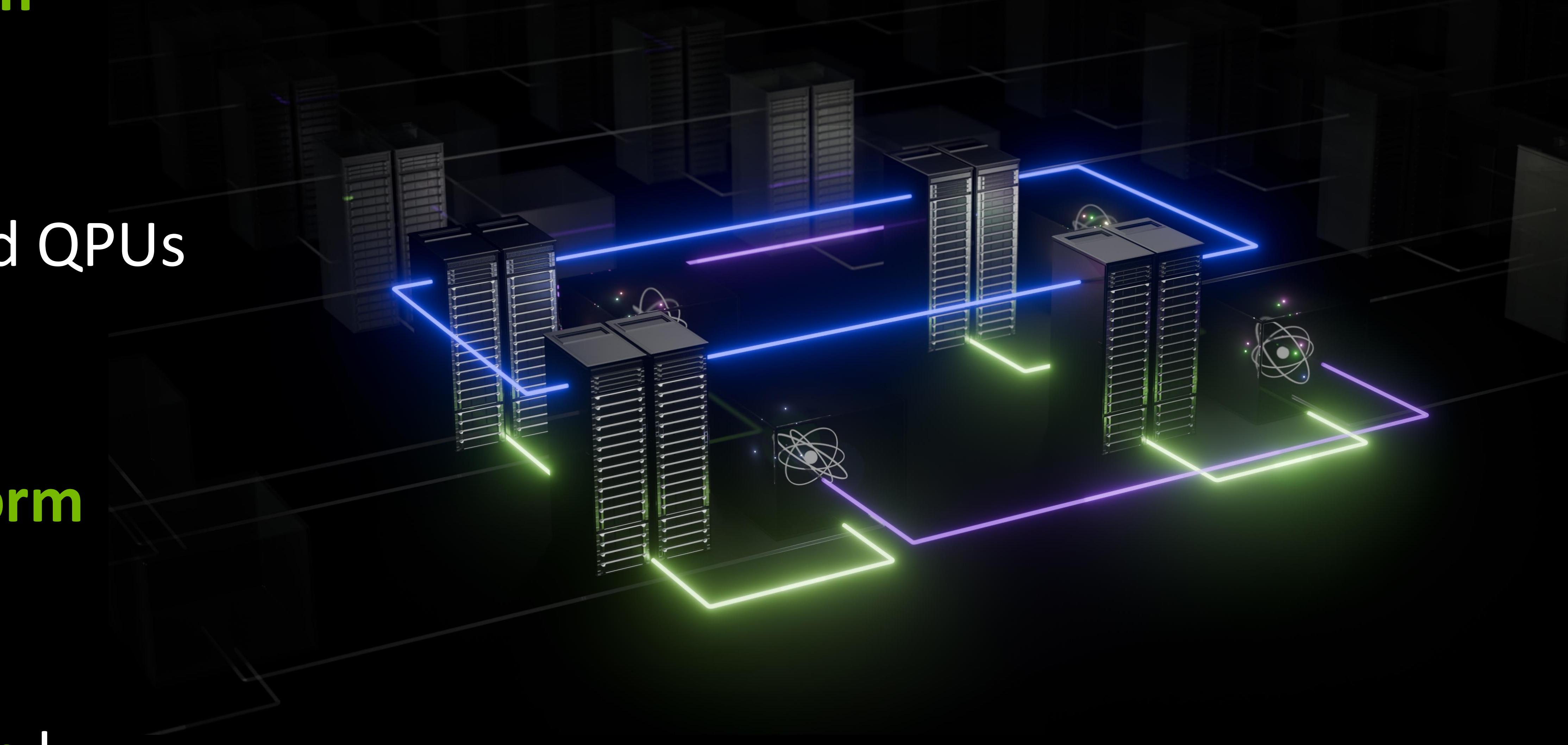
The Accelerated Quantum Supercomputer



# DGX SuperPOD Integrates the First Useful QPUs

The Accelerated Quantum Supercomputer

- DGX SuperPODs **connect quantum hardware as a coprocessor**
- **Hybrid algorithms** need GPUs and QPUs (Quantum Processing Units)
- NVIDIA's **CUDA-Q software platform** connects applications seamlessly
- NVIDIA solutions **de-risk quantum** by being qubit-agnostic





# Agenda

- Accelerated Quantum Supercomputer
- CUDA-Q Platform
- Accelerating Quantum Simulations with GPUs
- Installing CUDA-Q and the NVIDIA Quantum Cloud
- CUDA-Q Resources

# NVIDIA Quantum

Enabling the quantum ecosystem

QPU vendors

App developers

Academics

...

## CUDA-Q

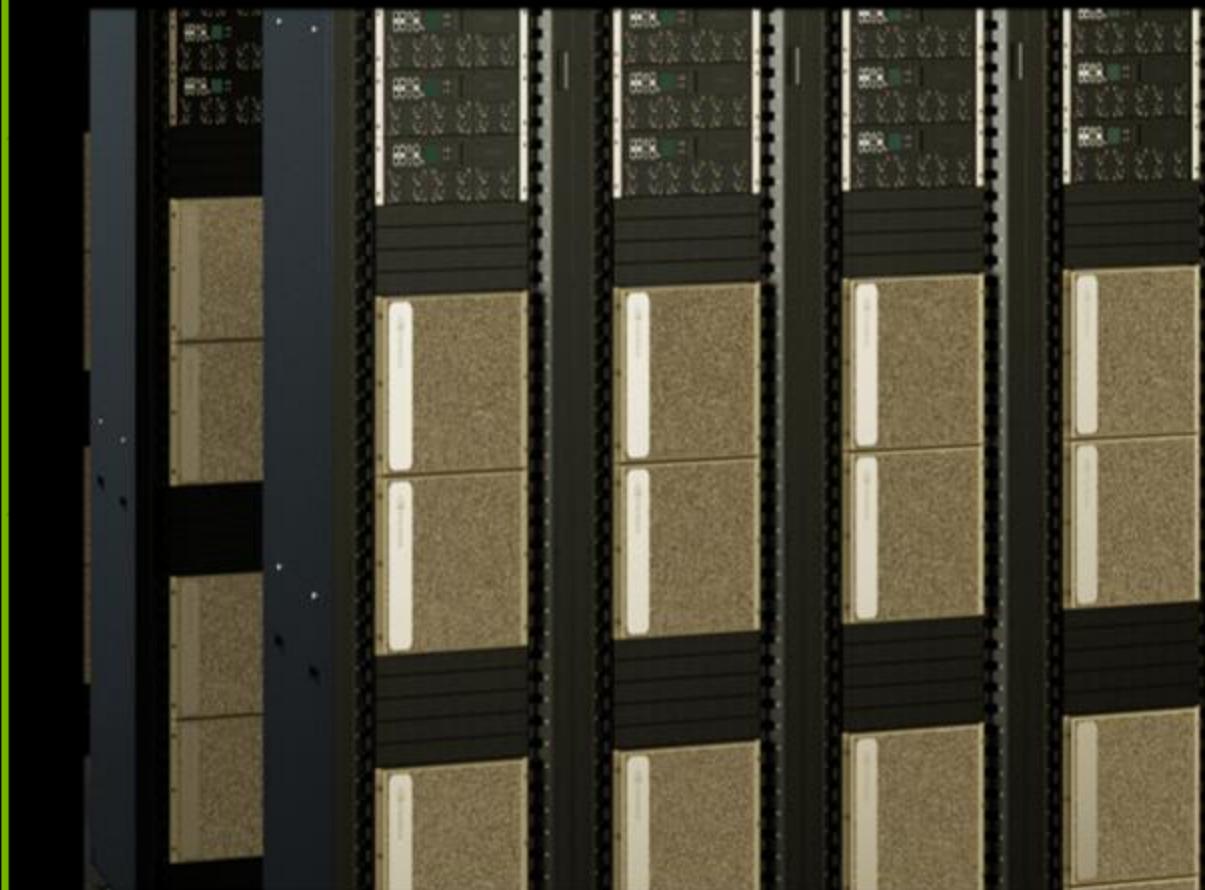
Libraries

Programming Model

Tools

Infrastructure

cuQuantum



Quantum Simulation

DGX-Q



Quantum Integrated  
Computing

CUDA-X

World's fastest AI and HPC libraries



# Agenda

- Accelerated Quantum Supercomputer
- CUDA-Q Platform
- Accelerating Quantum Simulations with GPUs
- Installing CUDA-Q and the NVIDIA Quantum Cloud
- CUDA-Q Resources

Github Repository with tutorial notebooks:  
<https://github.com/NERSC/sc24-quantum-tutorial>

# Onboarding Perlmutter

# Building your First CUDA-Q Kernel

```
import cudaq

qubit_count = 2

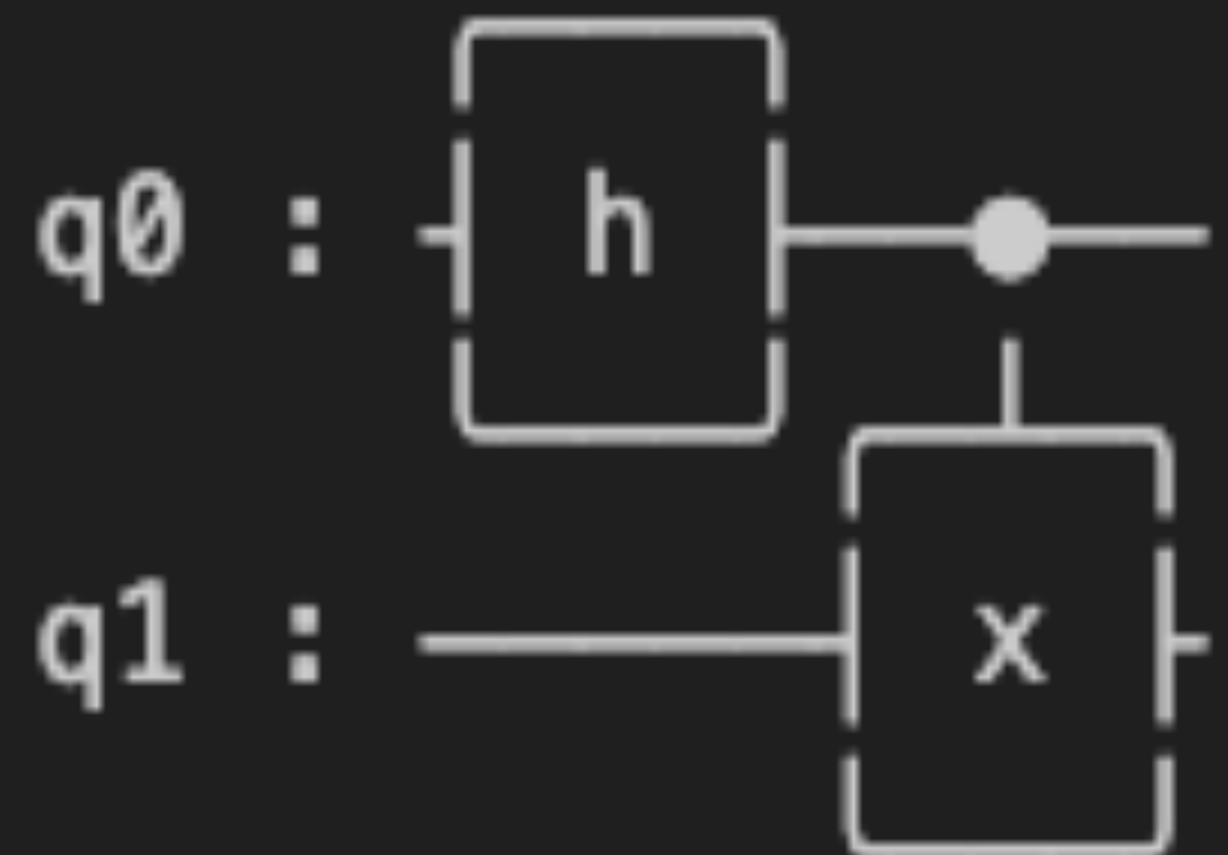
# Define the kernel
@cudaq.kernel
def my_kernel(qubit_count: int):
    # Allocate our `qubit_count` to the kernel.
    qubits = cudaq.qvector(qubit_count)

    # Apply a Hadamard gate to the qubit indexed by 0.
    h(qubits[0])

    # Apply a Controlled-X gate between qubit 0 (acting
    # as the control) and each of the remaining qubits.
    for i in range(qubit_count - 1):
        x.ctrl(qubits[i], qubits[i + 1])

    # Measure the qubits
    # If we don't specify measurements, all qubits are measured in the Z-basis by default.
    mz(qubits)

print(cudaq.draw(my_kernel, qubit_count))
```



# Sampling your First CUDA-Q Kernel

```
# First set the backend for kernel execution
cudaq.set_target('qpp-cpu') # selects a CPU backend

if cudaq.num_available_gpus() > 0:
    cudaq.set_target(`nvidia') # selects a GPU backend

# cudaq.set_target('nvqc') # selects the NVIDIA Quantum Cloud
# cudaq.set_target('ionq') # select an available QPU backend

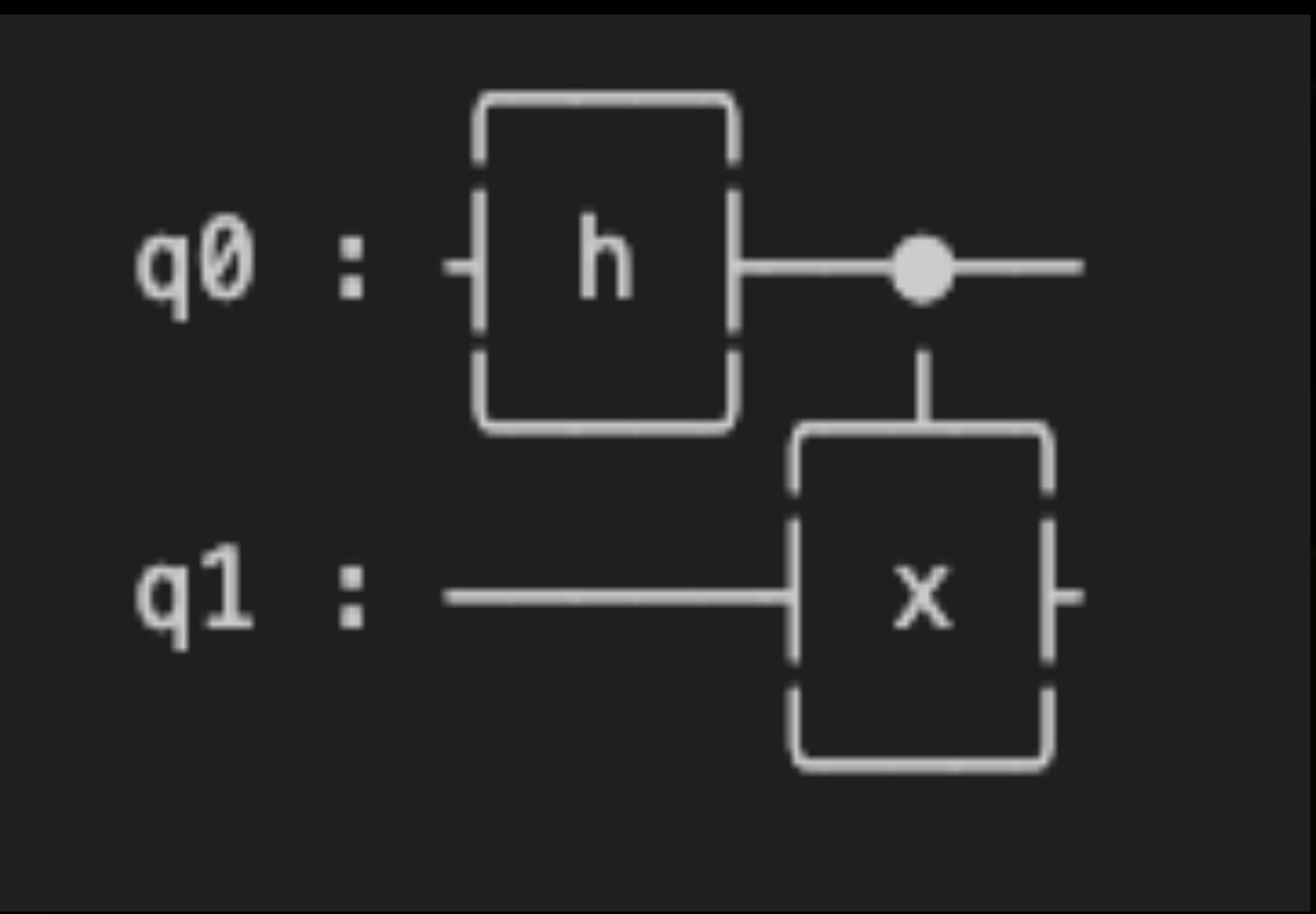
qubit_count = 2

results = cudaq.sample(my_kernel, qubit_count, shots_count = 10000)

print(results) # Example: {00:5005, 11: 4995}

print(results.most_probable()) # prints: `00`

print(results.probability(results.most_probable())) # prints: `0.5005`
```



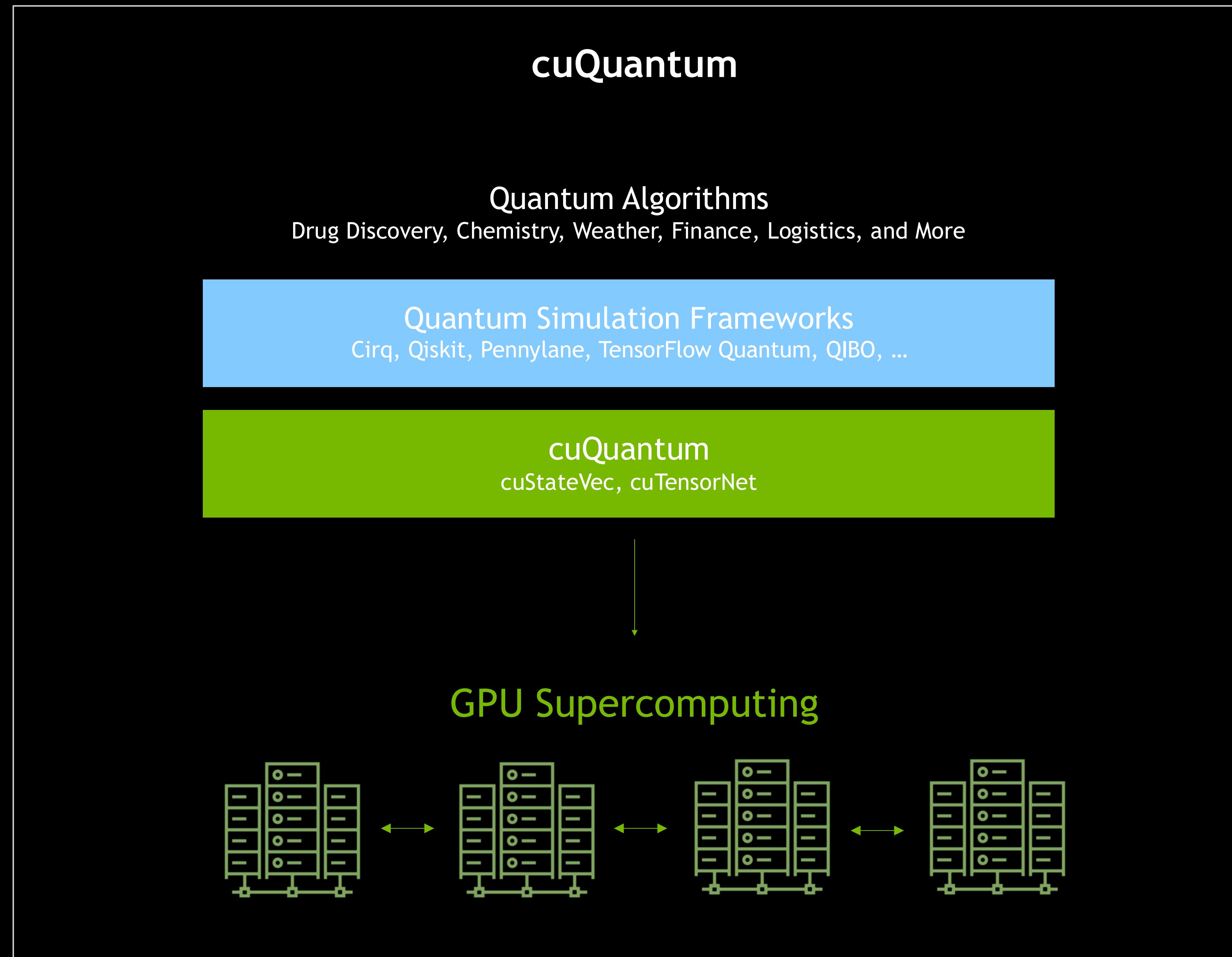
# Accelerating Quantum Simulations

Jupyter Notebook on Github

<https://github.com/NERSC/sc24-quantum-tutorial>

# cuQuantum

Research the Quantum Computer of Tomorrow on the most Powerful Computer Today



## CUQUANTUM FEATURES

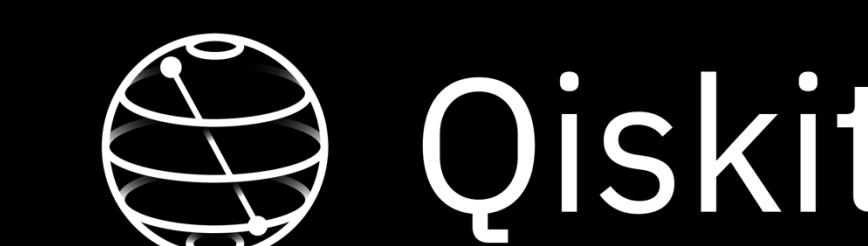
SDK for GPU Accelerated Quantum Simulation

Simulate Ideal or Noisy Qubits with State Vector or Tensor Network methods

Supports GPU Supercomputing with Multi-Node Multi-GPU Circuit Simulation

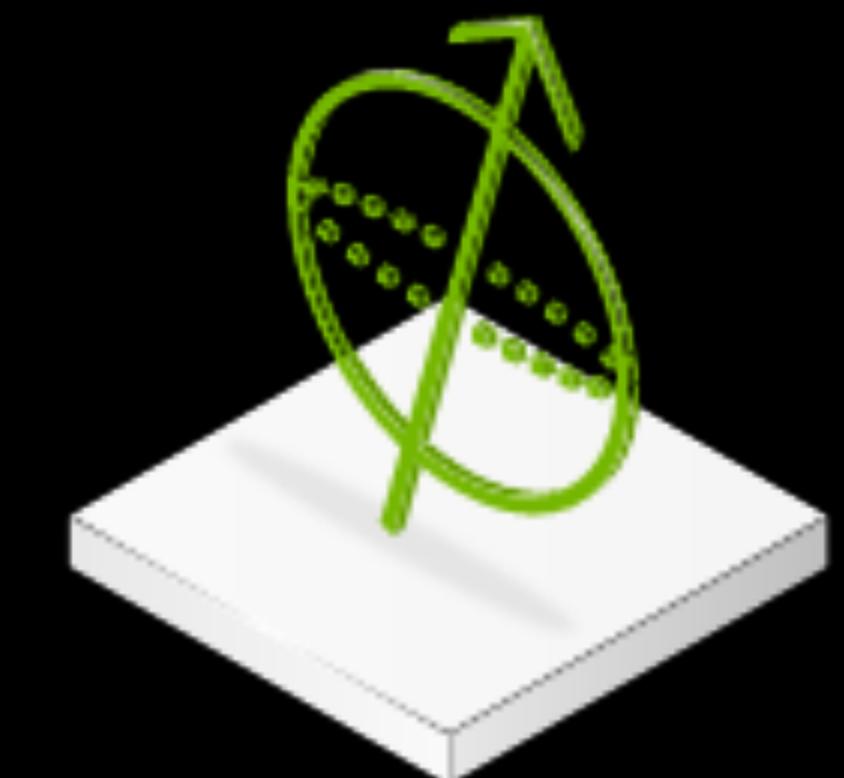
Accelerates all leading frameworks

Optimized frameworks in cuQuantum Appliance:  
[catalog.ngc.nvidia.com/orgs/nvidia/containers/cuquantum-appliance](https://catalog.ngc.nvidia.com/orgs/nvidia/containers/cuquantum-appliance)



# cuStateVec

Part of NVIDIA cuQuantum SDK



- cuStateVec: a library to accelerate statevector-based quantum circuit simulation

- Most computations are “in-place” to reduce memory usage
- Provides low-level primitives to cover common use cases:

- 1) **Apply gate matrix**
- 2) Apply diagonal/general permutation matrix
- 3) Apply exponential of Pauli matrix product
- 4) Expectation using matrix as observable
- 5) Expectation on Pauli basis
- 6) Sampling
- 7) Measurement on a Z-product basis
- 8) Batched single qubit measurement
- 9) State vector segment extraction/update
- 10) **Qubit reordering on single/multiple device(s)**

## C API

```
custatevecStatus_t custatevecApplyMatrix(  
    custatevecHandle_t handle,  
    void *sv,  
    cudaDataType_t svDataType,  
    const uint32_t nIndexBits,  
    const void *matrix,  
    cudaDataType_t matrixDataType,  
    custatevecMatrixLayout_t layout,  
    const int32_t adjoint,  
    const int32_t *targets,  
    const uint32_t nTargets,  
    const int32_t *controls,  
    const int32_t *controlBitValues,  
    const uint32_t nControls,  
    custatevecComputeType_t computeType,  
    void *extraWorkspace,  
    size_t extraWorkspaceSizeInBytes)
```

## Python API

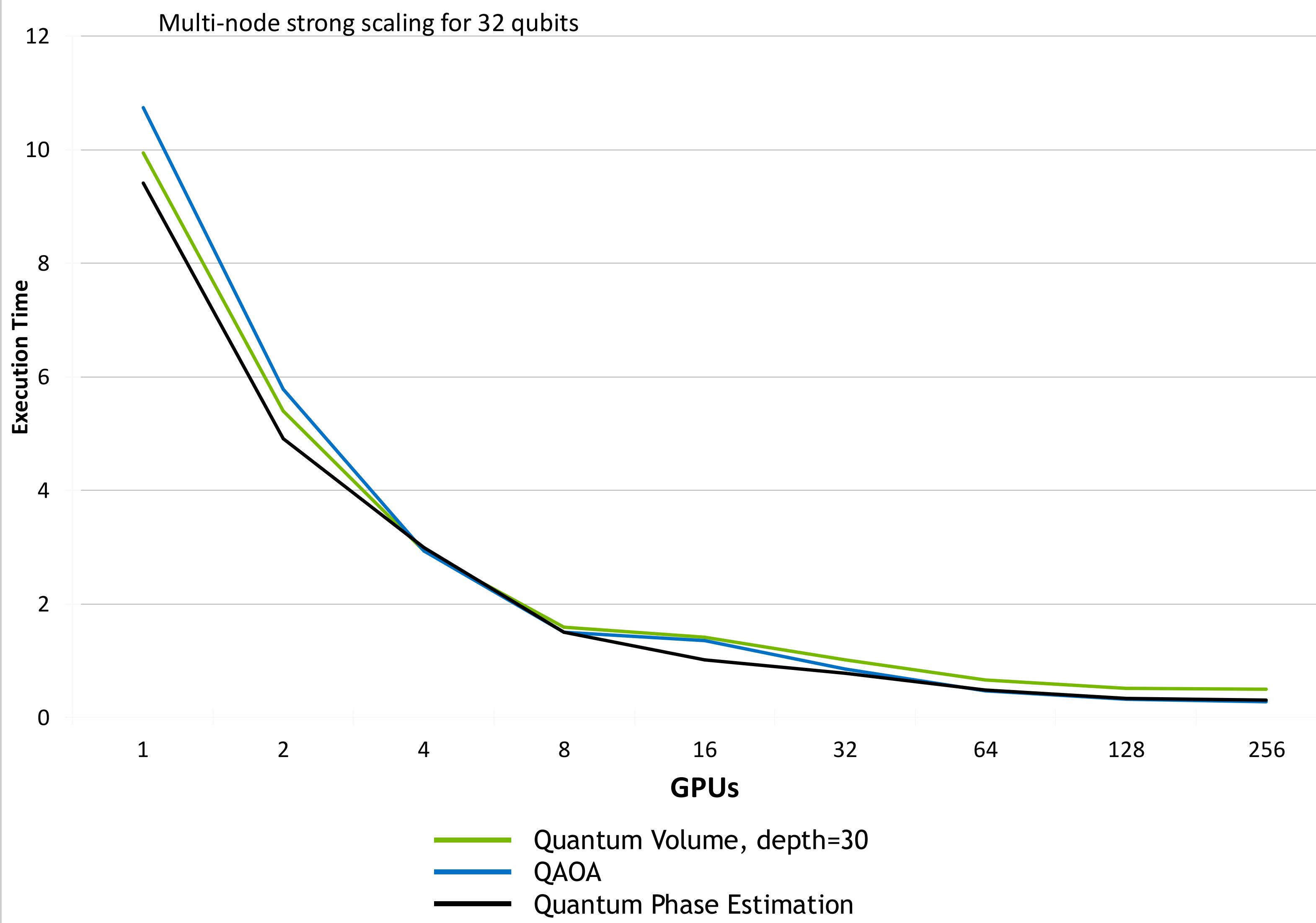
```
cuquantum.custatevec.apply_matrix(  
    handle,  
    sv,  
    sv_data_type,  
    n_index_bits,  
    matrix,  
    matrix_data_type,  
    layout,  
    adjoint,  
    targets,  
    n_targets,  
    controls,  
    control_bit_values,  
    n_controls,  
    compute_type,  
    workspace,  
    workspace_size)
```

- Easy integration & adoption for a wide variety of frameworks & programming languages
- Also available in the cuQuantum Appliance container (standalone & Cirq/Qsim backend)

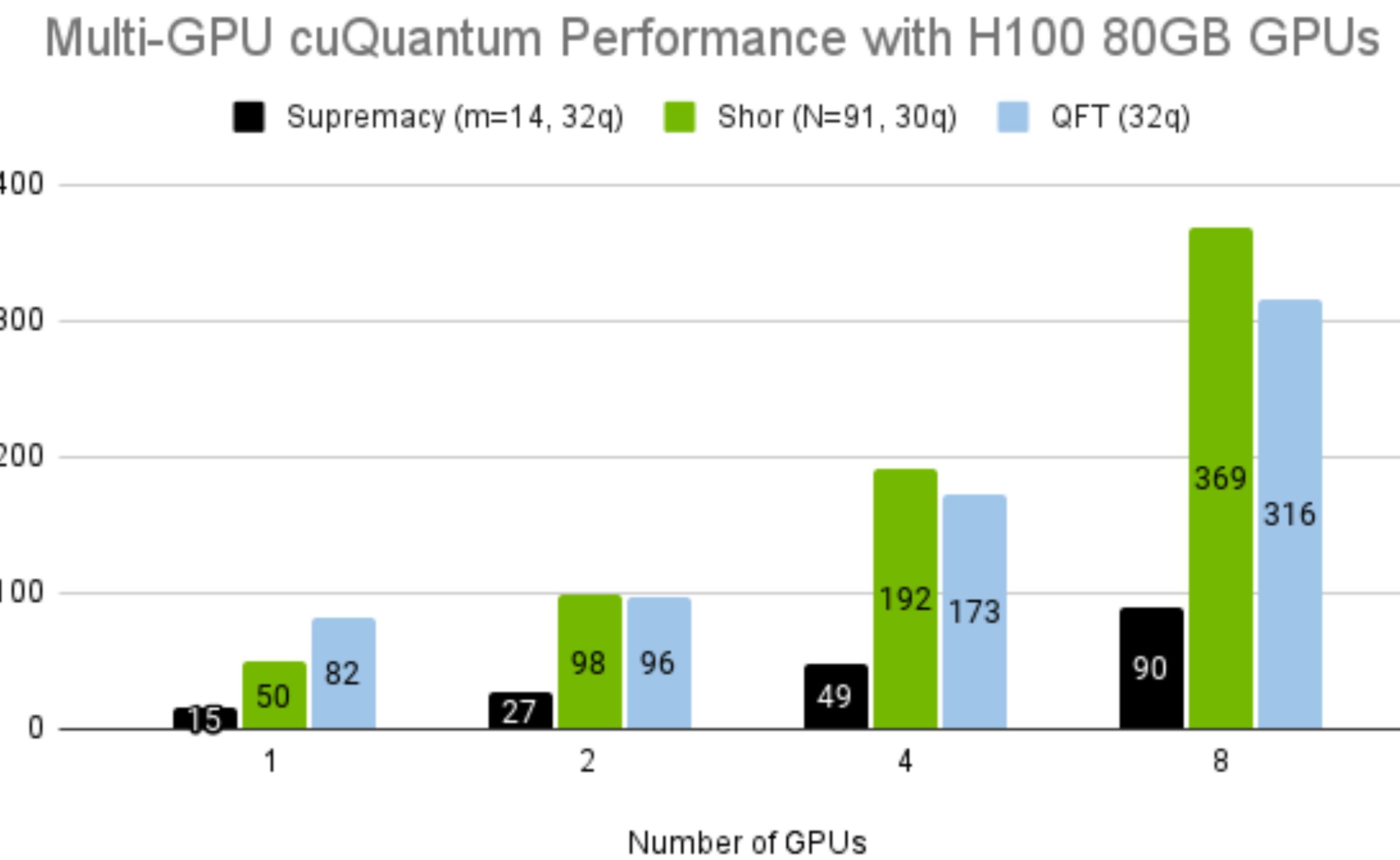
# cuStatevec

Research the Quantum Computer of Tomorrow on the most Powerful Computer Today

## World Class Performance with Now with Multi-Node Multi-GPU



## World Class Performance

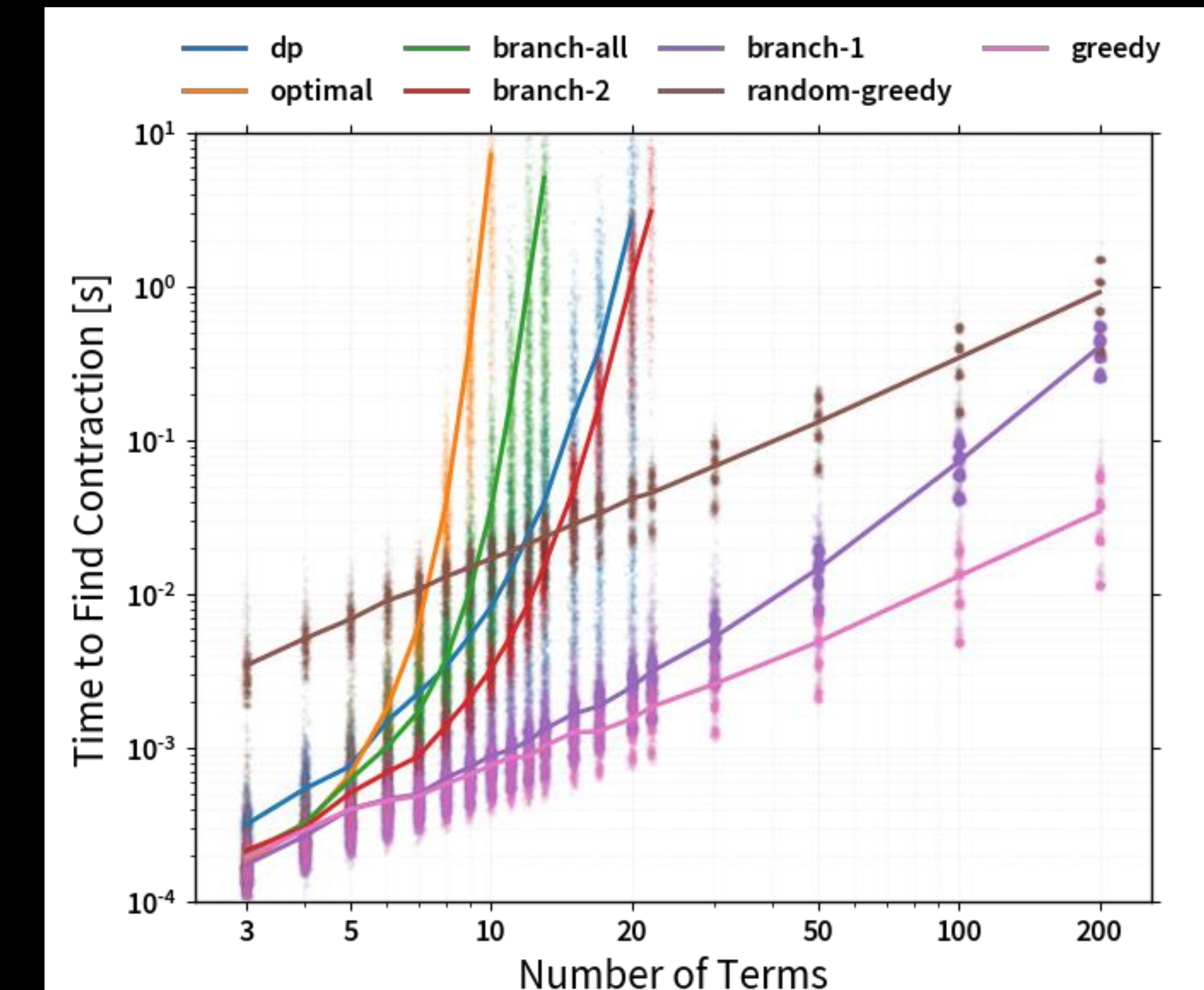
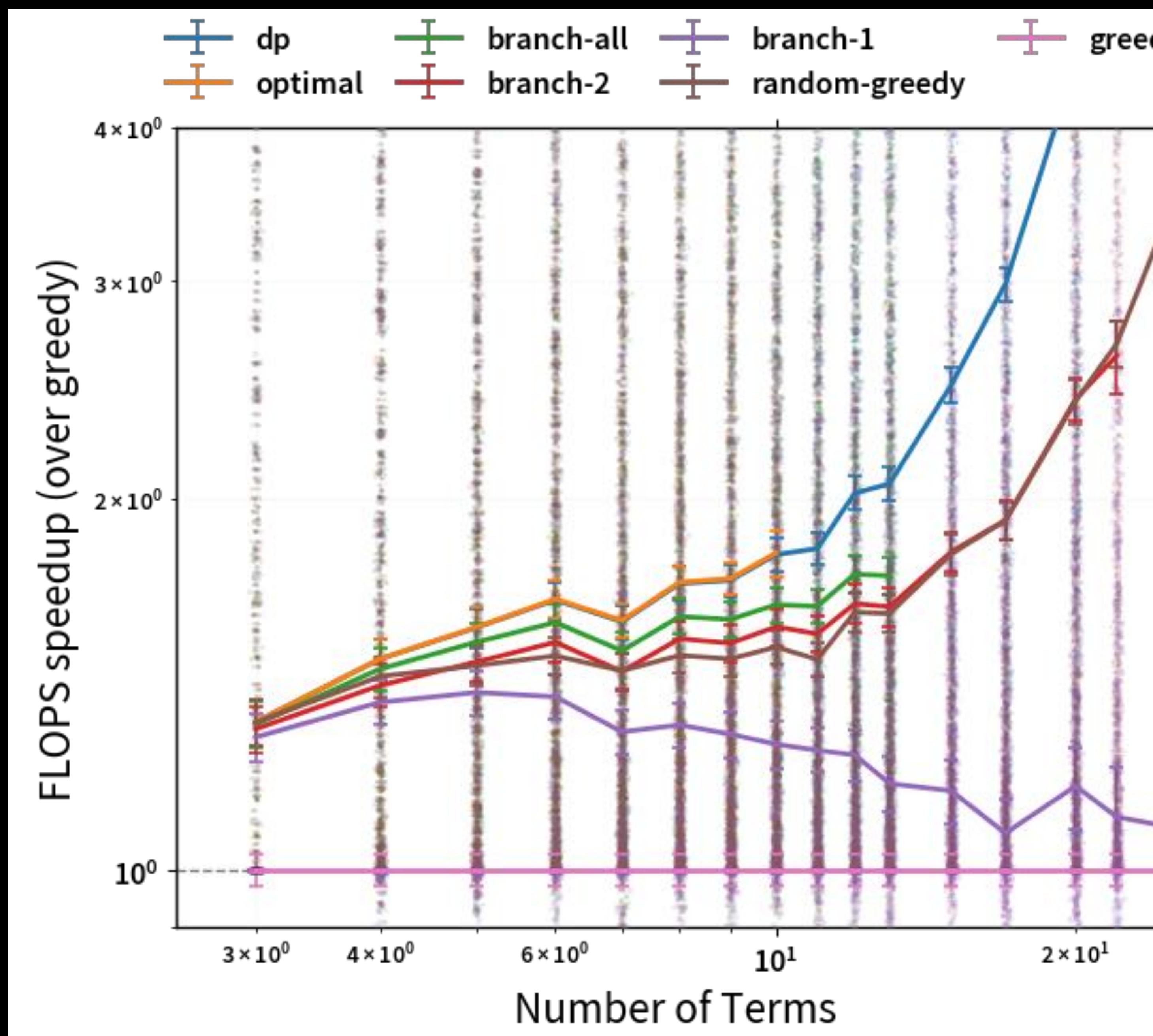


# cuTensorNet

A Library to accelerate tensor network based quantum circuit simulation

## Motivation: What is contraction path optimization and why is it needed ?

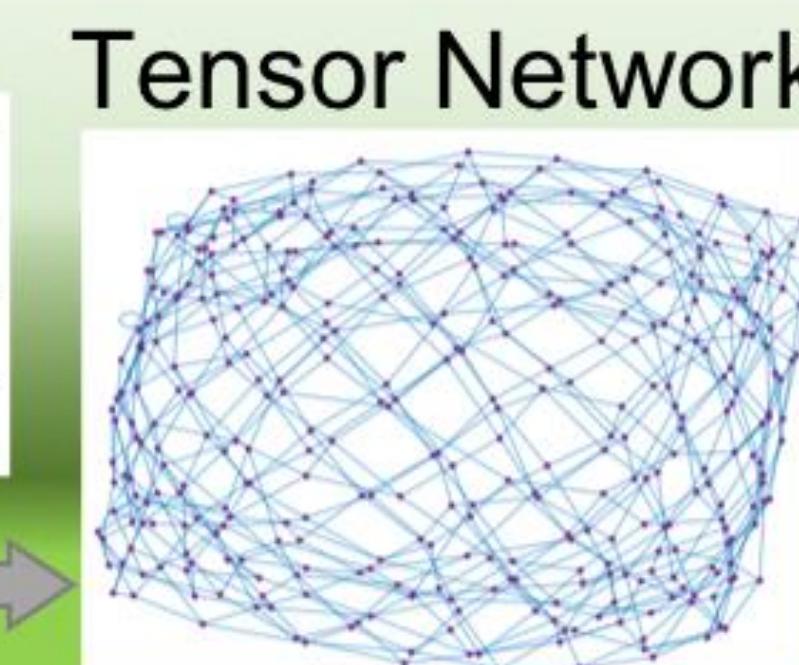
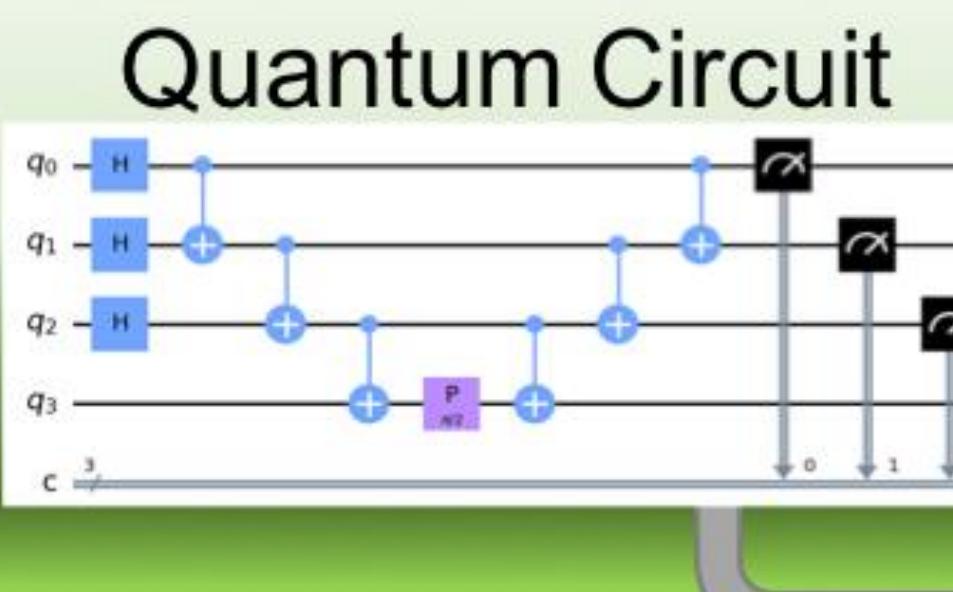
- Reordering the contraction is important to minimize flops as well as memory requirements.
- Finding the optimal path is NP hard problem, which is why they only can be used for small number of tensors



# cuTensorNet

A Library to accelerate tensor network based quantum circuit simulation

Quantum Algorithm



## cuTensorNet



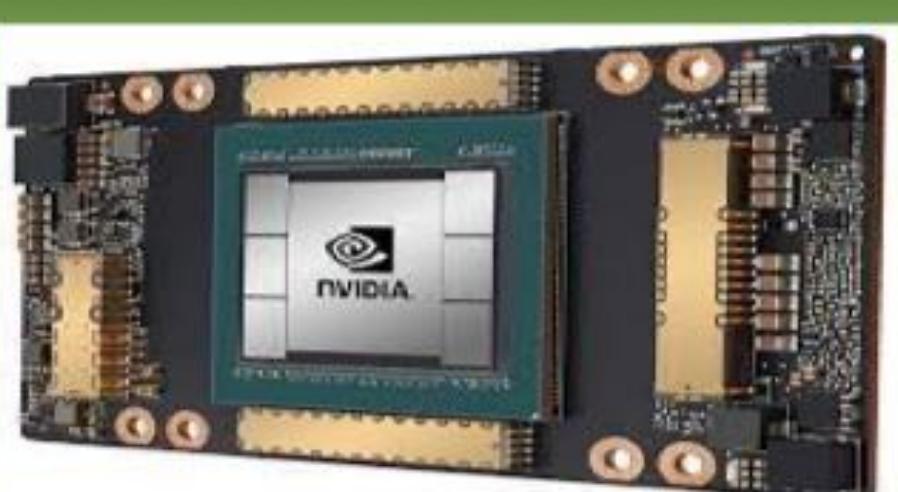
Algorithmic  
advancements



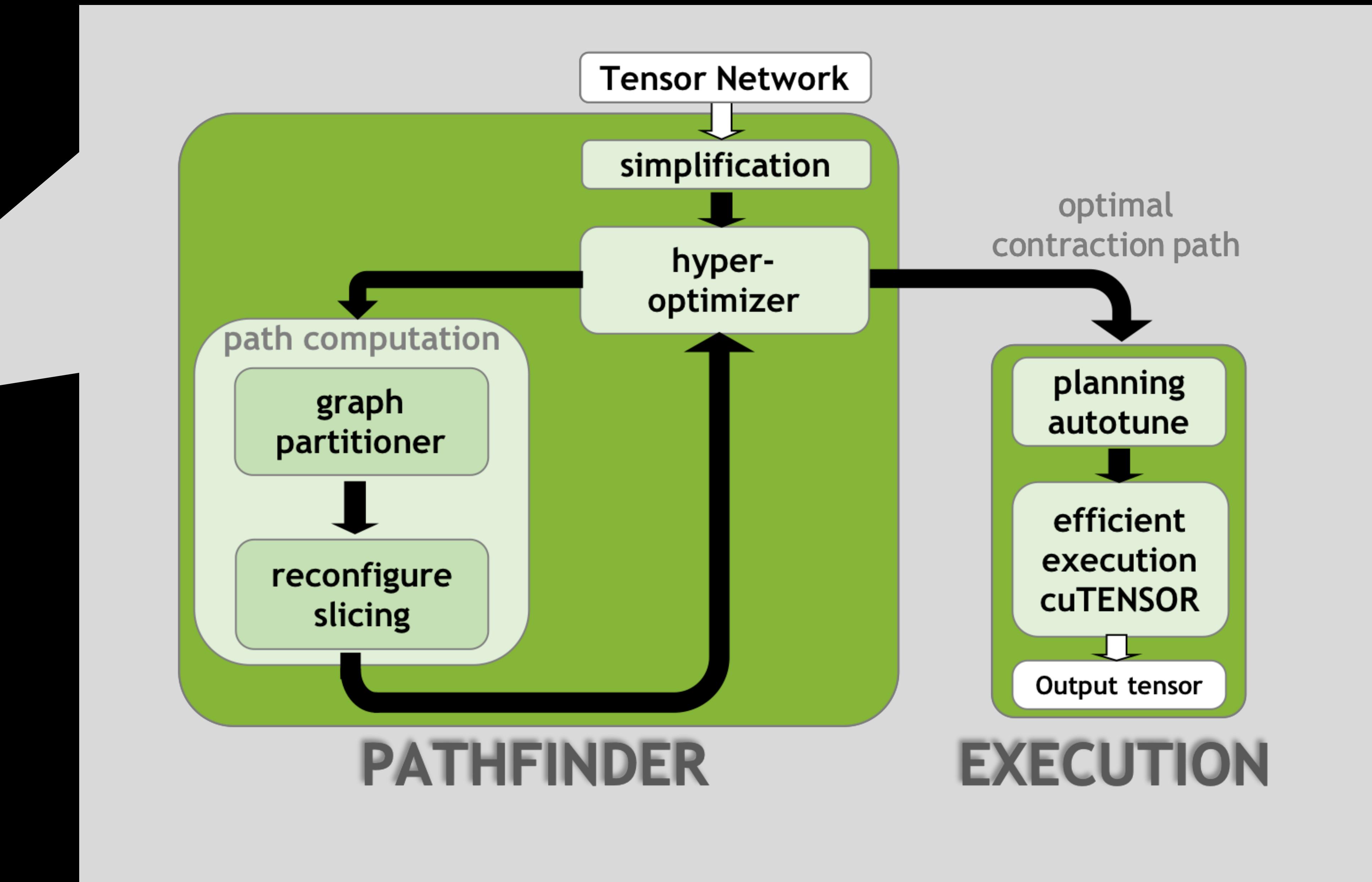
Optimal  
contraction

Kernels  
optimization

NVIDIA GPU

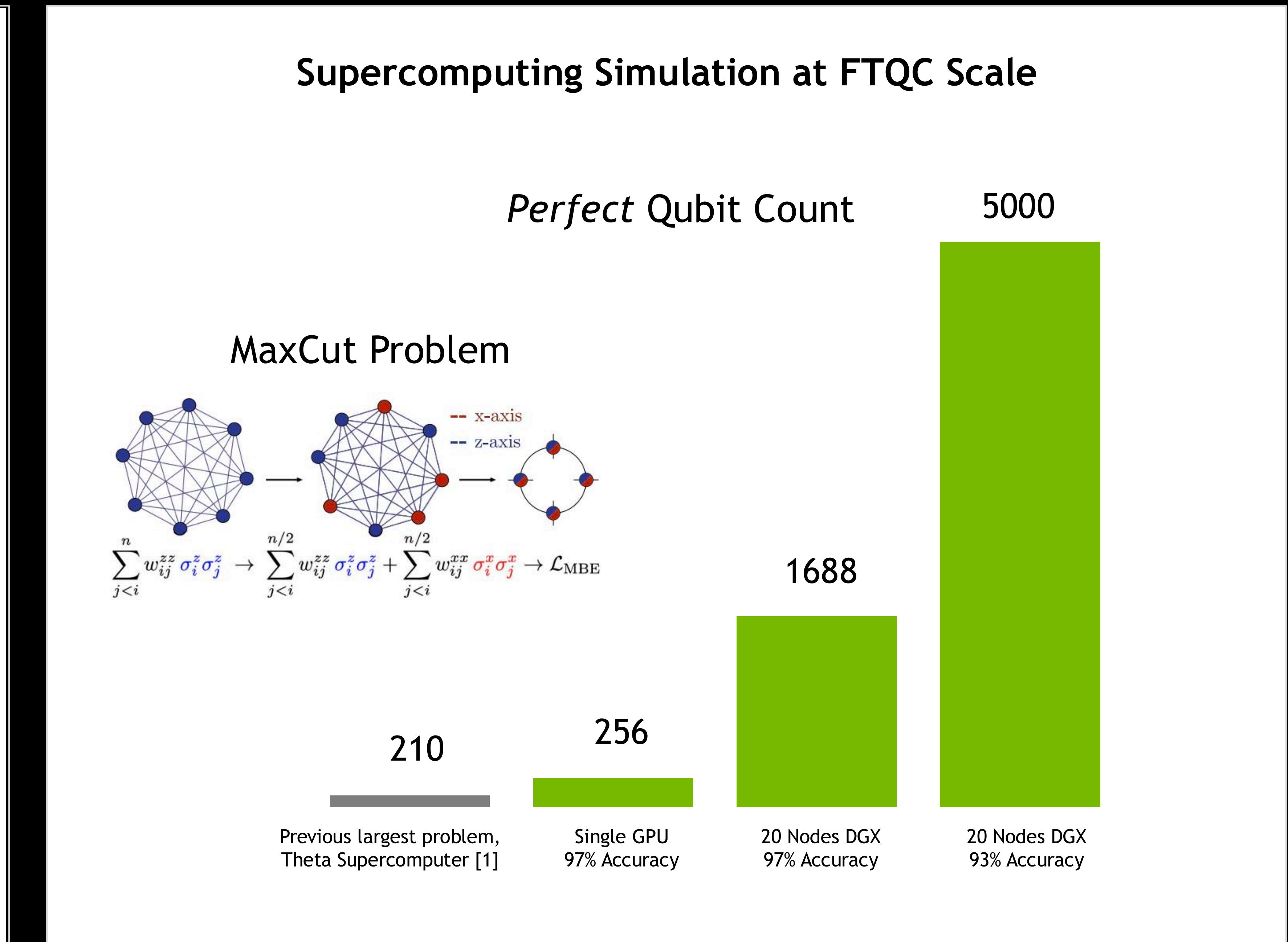
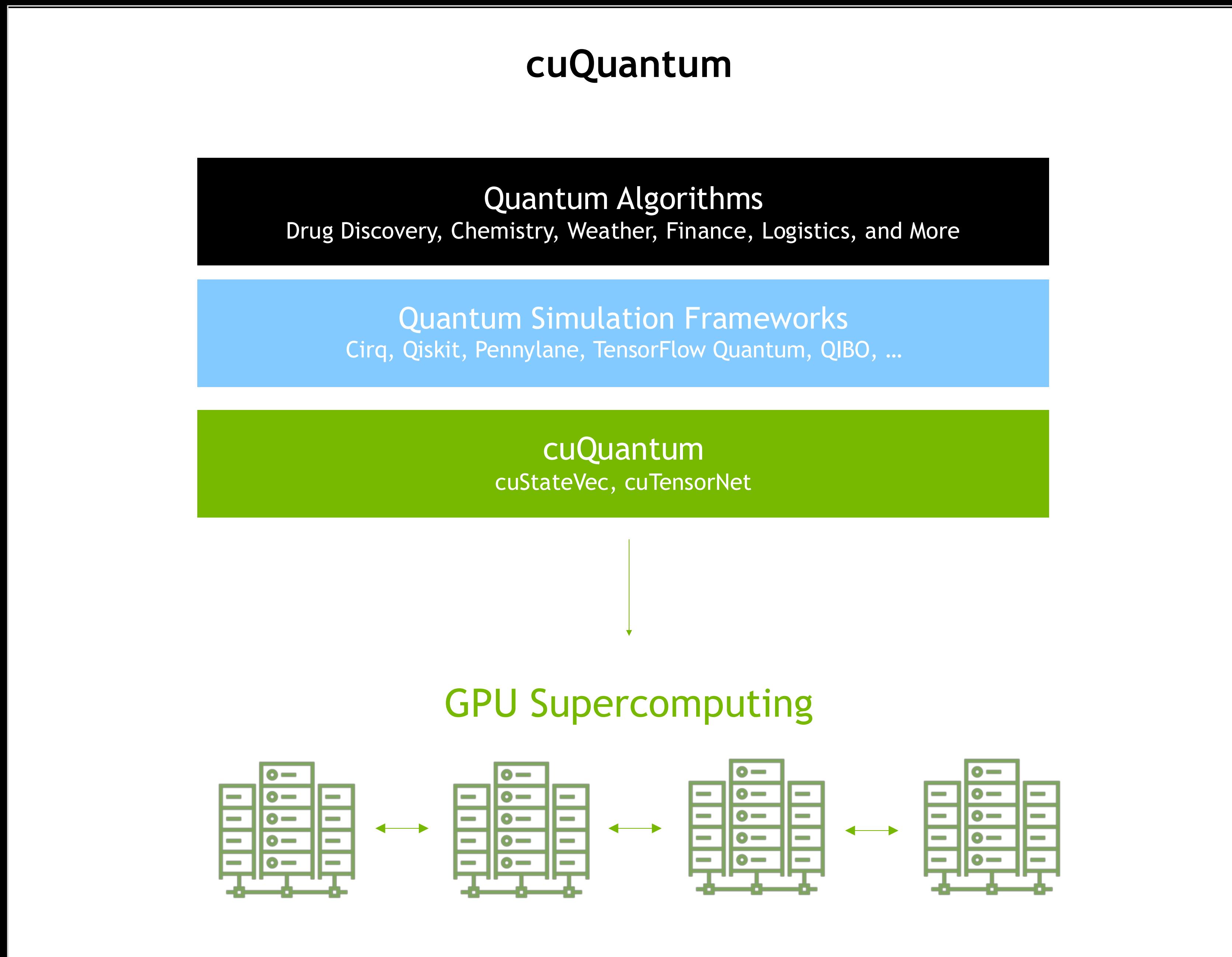


Supercomputer



# cuTensorNet

Research the Quantum Computer of Tomorrow on the most Powerful Computer Today

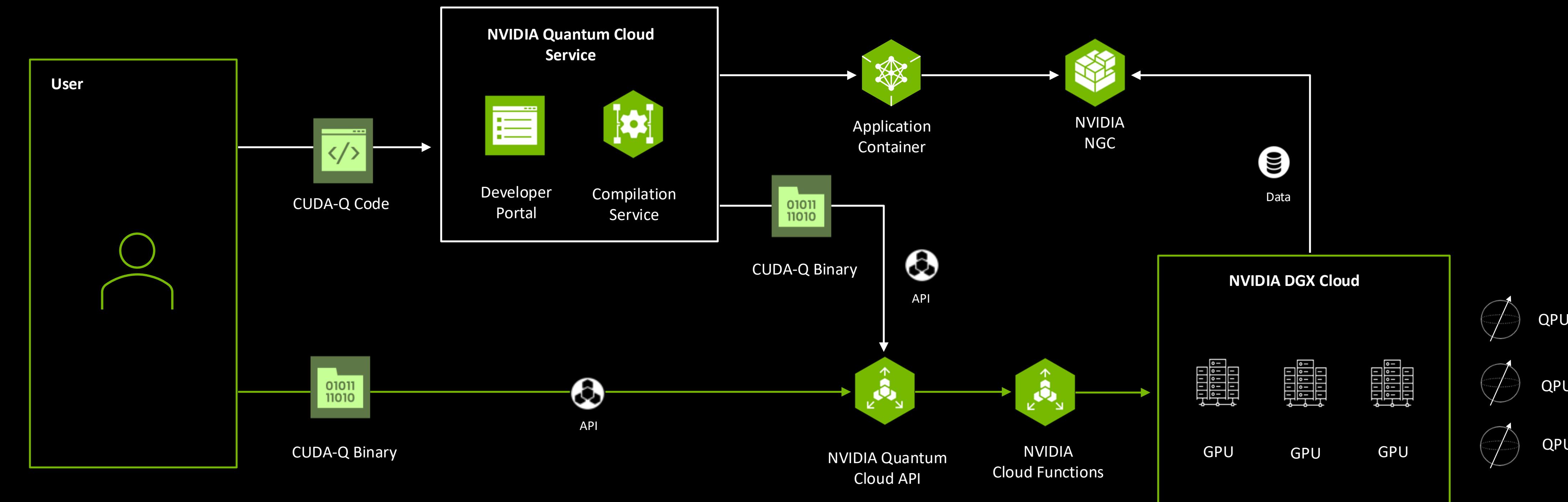


[1] Danylo Lykov et al, Tensor Network Quantum Simulator With Step-Dependent Parallelization, 2020  
<https://arxiv.org/pdf/2012.02430.pdf>

[2] Taylor Patti et al, Variational Quantum Optimization with Multibasis Encodings, 2022  
<https://arxiv.org/abs/2106.13304>

# NVIDIA Quantum Cloud

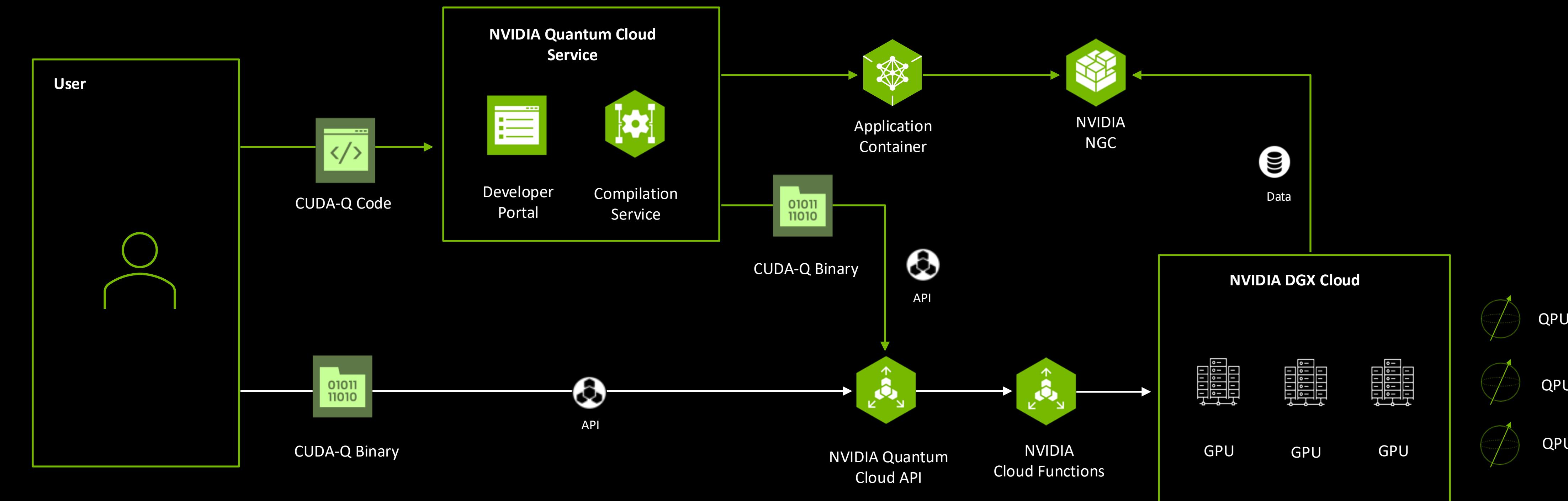
Seamless Access to the most Powerful Quantum Resources



- Access to the most powerful Quantum Resources
- Develop locally, run any CUDA-Q app seamlessly in Cloud
- Call pre-built Quantum Cloud APIs from your application
  - QPU Partner Integrations coming soon
  - Quantum Cloud Web Developer Portal coming soon
- Run workloads on GPU Supercomputers
- Integrated ISV applications
- EA Available Today! Apply for Early Access at [developer.nvidia.com/quantum-cloud-early-access-join](https://developer.nvidia.com/quantum-cloud-early-access-join)

# NVIDIA Quantum Cloud

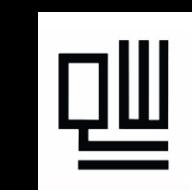
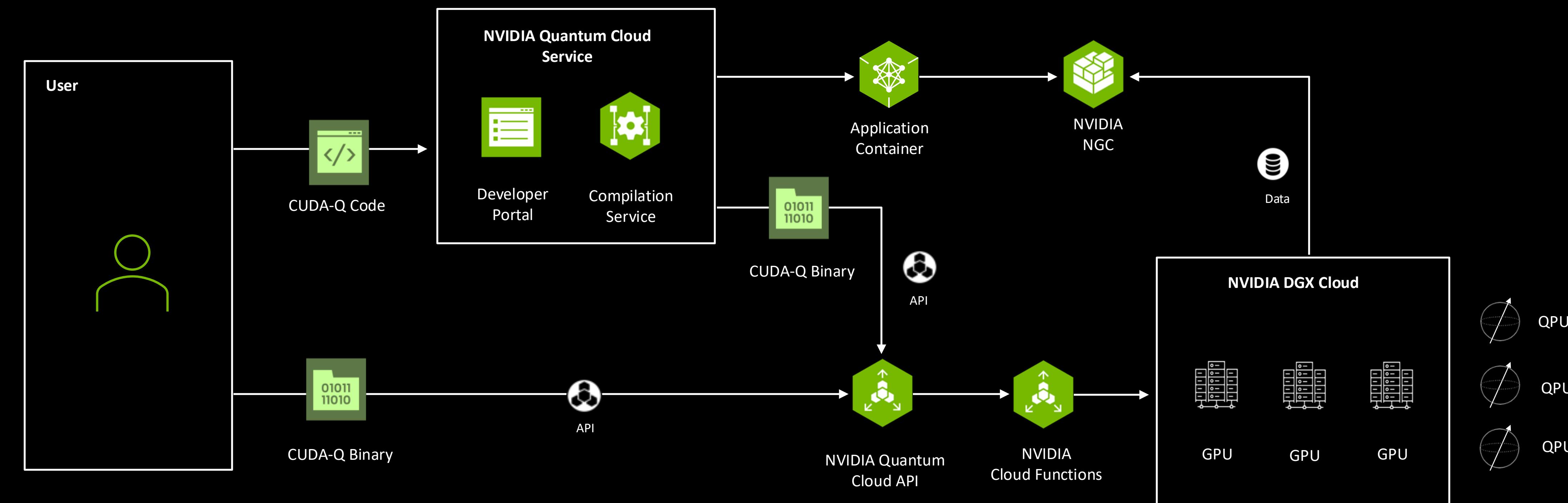
Seamless Access to the most Powerful Quantum Resources



- Access to the most powerful Quantum Resources
- Develop locally, run any CUDA-Q app seamlessly in Cloud
- Call pre-built Quantum Cloud APIs from your application
  - QPU Partner Integrations coming soon
  - Quantum Cloud Web Developer Portal coming soon
- Run workloads on GPU Supercomputers
- Integrated ISV applications
- EA Available Today! Apply for Early Access at [developer.nvidia.com/quantum-cloud-early-access-join](https://developer.nvidia.com/quantum-cloud-early-access-join)

# NVIDIA Quantum Cloud

Seamless Access to the most Powerful Quantum Resources



Find out more

## NVIDIA Quantum

<https://www.nvidia.com/en-us/solutions/quantum-computing/>

## CUDA-Q Academic

Educational resources for CUDA-Q  
(<https://github.com/NVIDIA/cuda-q-academic>)

## NVIDIA Quantum Cloud

**Early Access** - <https://www.nvidia.com/en-us/solutions/quantum-computing/cloud/>





# Agenda

## Morning Sessions

- Introductions
- Overview of CUDA-Q
- Onboarding to Perlmutter
- Accelerating Quantum Simulations Notebook
- NVIDIA Quantum Cloud and Local install of CUDA-Q

## Afternoon Sessions

- Intro to large scale clusters, how to navigate and use them
- Example: Quantum chemistry at NERSC
- Industry Use Case: Simulating Hamiltonians of molecules with 30,000 terms
- QuEra's Quantum Hardware
- Example: Quantum reservoir computing

# Outline for NERSC Section of the CUDA-Q Tutorial

- Description of NERSC systems
- How to access NERSC training accounts
- Physics and Chemistry Variational Quantum Eigensolver examples:

<https://github.com/NERSC/sc24-quantum-tutorial>

# NERSC Overview



QIS Tutorial at SC24  
Nov 17th 2024

Katie Klymko, Ermal Rrapaj  
Advanced Technologies Group  
NERSC

# National Energy Research Scientific Computing Center

- NERSC is a national supercomputer center funded by the U.S. Department of Energy Office of Science (SC)
  - Supports SC research mission
  - Part of Berkeley Lab
- Researchers with funding from SC who need supercomputing resources can use NERSC
  - Other researchers can apply if research is in SC mission
- NERSC supports 10,000 users, 1,000 projects
  - From all 50 states + international; 60% from universities
  - Hundreds of users log on each day

# NERSC User Demographics

~10,000 Annual Users from ~800 Institutions + National Labs



**32%**  
Graduate  
Students



**19%**  
Postdoctoral  
Fellows



**15%**  
Staff  
Scientists



**13%**  
University  
Faculty



**8%**  
Undergraduate  
Students



**5%**  
Professional  
Staff



**60%**  
Universities



**29%**  
DOE Labs



**5%**  
Other  
Government Labs



**4%**  
Industry

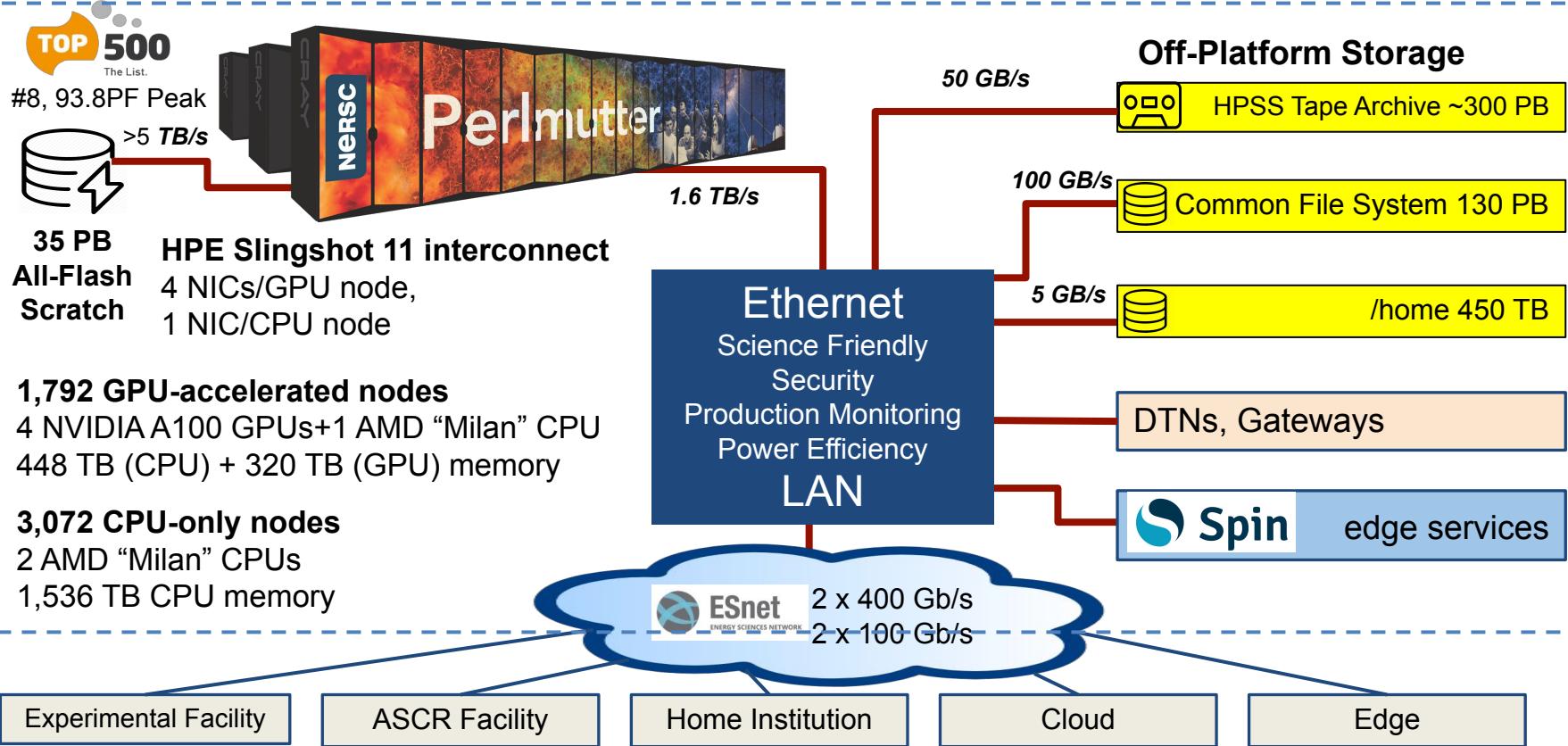


**1%**  
Small  
Businesses



**<1%**  
Private Labs

# NERSC Systems Ecosystem



# Perlmutter system configuration

## NVIDIA "Ampere" GPU Nodes

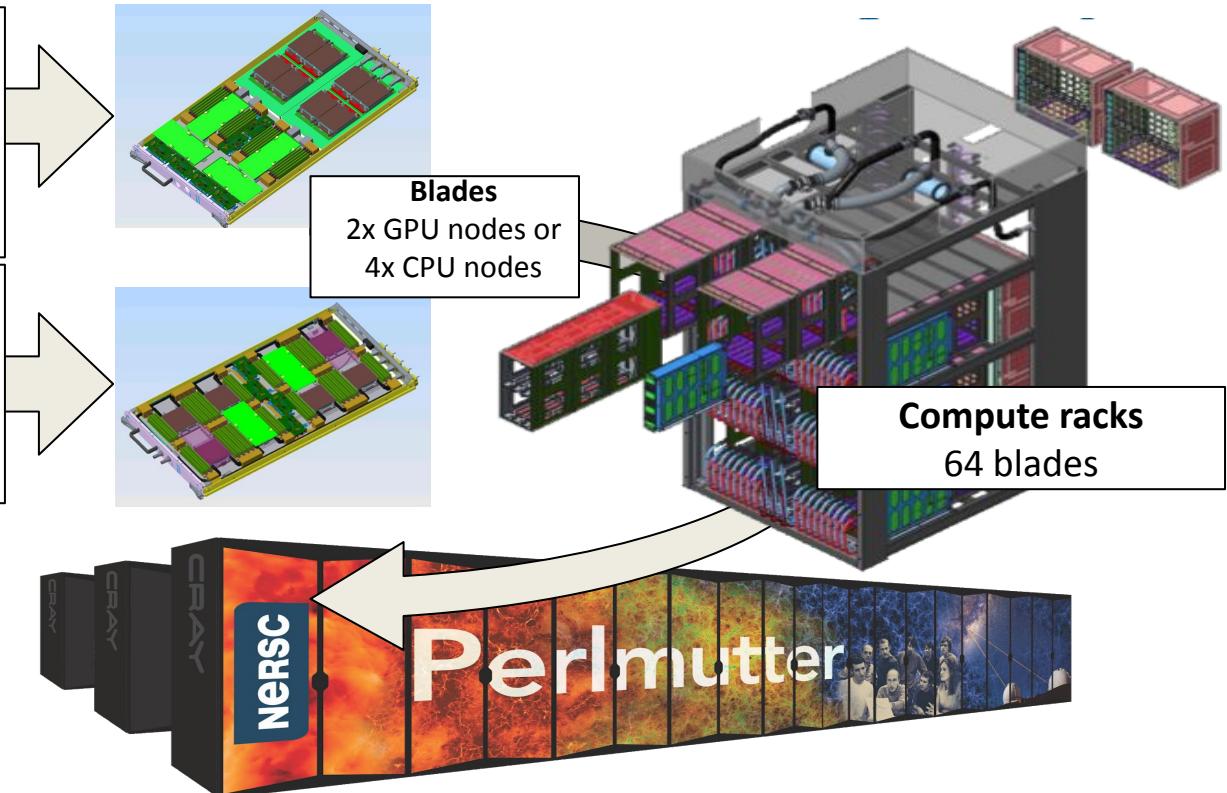
4x GPU + 1x CPU  
40 GiB HBM + 256 GiB DDR  
4x 200G "Slingshot" NICs

## AMD "Milan" CPU Node

2x CPUs  
> 256 GiB DDR4  
1x 200G "Slingshot" NIC

Centers of Excellence  
Network  
Storage  
App. Readiness  
System SW

Perlmutter system  
GPU racks  
CPU racks  
~6 MW

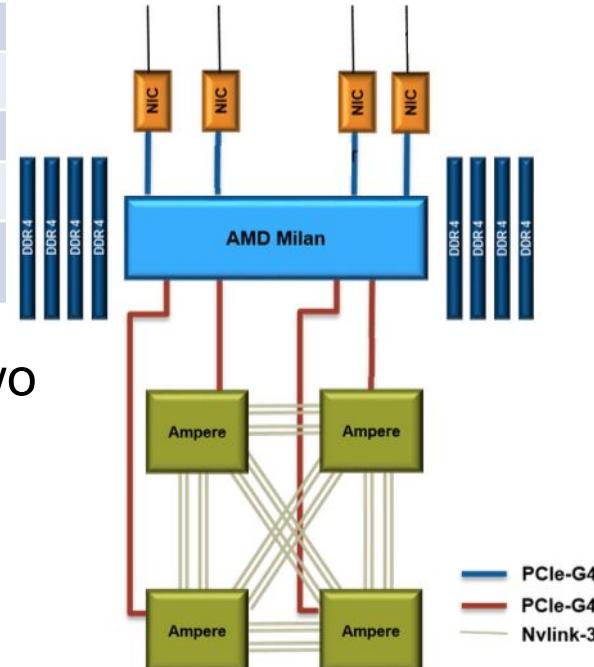


# Node Types

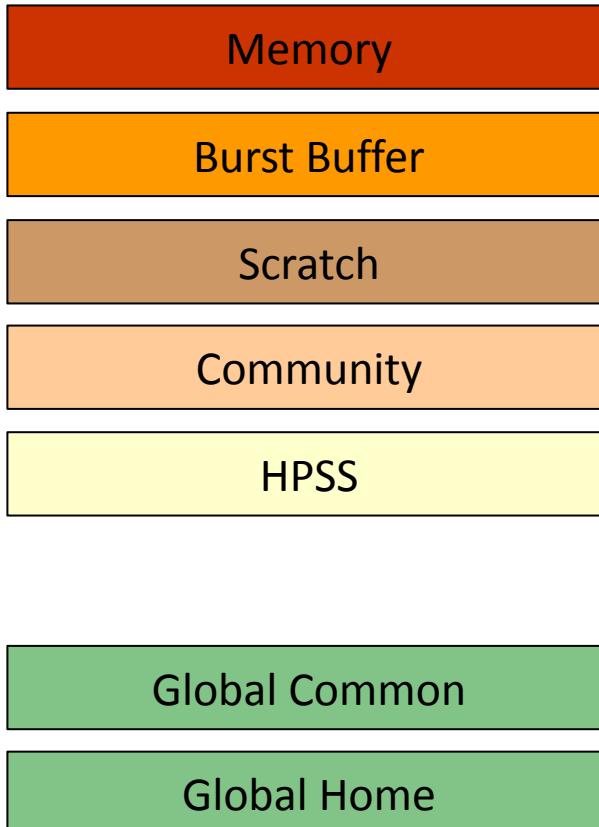
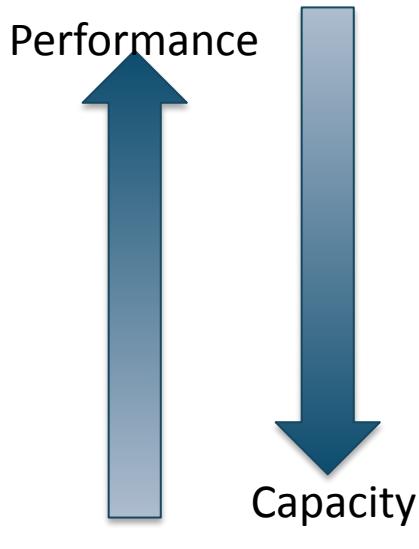
Partition	Nodes	CPU	RAM	GPU	NIC
GPU	1536	1x <a href="#">AMD EPYC 7763</a>	256GB	4x <a href="#">NVIDIA A100</a> (40GB)	4x <a href="#">HPE Slingshot 11</a>
	256	1x <a href="#">AMD EPYC 7763</a>	256GB	4x <a href="#">NVIDIA A100</a> (80GB)	4x <a href="#">HPE Slingshot 11</a>
CPU	3072	2x <a href="#">AMD EPYC 7763</a>	512GB	–	1x <a href="#">HPE Slingshot 11</a>
Login	40	1x <a href="#">AMD EPYC 7713</a>	512GB	4x <a href="#">NVIDIA A100</a> (40GB)	–
Large Memory	4	1x <a href="#">AMD EPYC 7713</a>	1TB	4x <a href="#">NVIDIA A100</a> (40GB)	1x <a href="#">HPE Slingshot 11</a>

- Each Milan CPU has 64 2.45 GHz cores with two hardware threads
- Access is managed via Slurm partitions
- Login nodes are GPU-enabled

GPU Node Architecture:



# Simplified NERSC File Systems



## 35 PB Flash Scratch

Lustre >5 TB/s

temporarily (purge)

**xGB (x=Total RAM) In-Memory Burst Buffer**  
/tmp RamFS

## 157 PB HDD Community

Spectrum Scale (GPFS)

150 GB/s, permanent

## 150 PB Tape Archive

HPSS Forever

## 20 TB SSD Software

Spectrum Scale

Permanent

Faster compiling / Source Code

# Global File Systems

## Global Home

- Permanent, relatively small storage
- Mounted on all platforms
- NOT tuned to perform well for parallel jobs
- Quota cannot be changed
- Snapshot backups (7-day history)
- **Perfect for storing data such as source code, shell scripts**
- **Addressed using \$HOME**

## Community File System (CFS)

- Permanent, larger storage
- Mounted on all platforms
- Medium performance for parallel jobs
- Quota can be changed
- Snapshot backups (7-day history)
- **Perfect for sharing data within research group**
- **Addressed using \$CFS**

# Local File Systems

## Scratch

- Large, temporary storage
- Optimized for read/write operations, NOT storage
- Not backed up
- Purge policy (8 weeks)
- **Perfect for staging data and performing computations**
- **Addressed using \$SCRATCH**

## Burst Buffer

- Temporary storage
- High-performance in-memory file system
- **Perfect for getting good performance in I/O-constrained codes**
- **Not shared between nodes**

There are many different ways to access NERSC. To use our resources, you need to either:

1. Log into the login nodes and interact with Slurm (covered here)
2. Use services (eg. Jupyter) that expose web interfaces (covered later)
3. Interact via our REST API – called the “Superfacility API”  
(see: <https://docs.nersc.gov/services/sfapi/>)

# Connecting with SSH

- To access Perlmutter, use:

saul.nersc.gov

perlmutter.nersc.gov



```
blaschke@laptop:~$ ssh <user>@saul.nersc.gov  
blaschke@laptop:~$ ssh <user>@perlmutter.nersc.gov
```

- To be able to open GUI applications on the login nodes use: `ssh -Y`



```
blaschke@laptop:~$ ssh -Y <user>@saul.nersc.gov
```

# Connecting with SSH

After successfully logging on, you  
be greeted by the terms of use, and  
the command-line input prompt:

```
For past outages, see: https://my.nersc.gov/outagelog-cs.php  
blaschke@perlmutter:login17:~>
```

From here you can interact with  
perlmutter...



```
blaschke@laptop:~$ ssh blaschke@saul.nersc.gov
*****
NOTICE TO USERS

Lawrence Berkeley National Laboratory operates this computer system under
contract to the U.S. Department of Energy. This computer system is the
property of the United States Government and is for authorized use only.
Users (authorized or unauthorized) have no explicit or implicit
expectation of privacy.

Any or all uses of this system and all files on this system may be
intercepted, monitored, recorded, copied, audited, inspected, and disclosed
to authorized site, Department of Energy, and law enforcement personnel,
as well as authorized officials of other agencies, both domestic and foreign.
By using this system, the user consents to such interception, monitoring,
recording, copying, auditing, inspection, and disclosure at the discretion
of authorized site or Department of Energy personnel.

Unauthorized or improper use of this system may result in administrative
disciplinary action and civil and criminal penalties. By continuing to use
this system you indicate your awareness of and consent to these terms and
conditions of use. LOG OFF IMMEDIATELY if you do not agree to the conditions
stated in this warning.

*****
Login connection to host x3114c0s15b0n0:

#####
# [REDACTED] #
#####

Hello, World!
Welcome to perlmutter!
#####
# [REDACTED] #
#####

For past planned outages, see: https://www.nersc.gov/live/status/motd/
For past outages, see: https://my.nersc.gov/outagelog-cs.php
blaschke@perlmutter:login17:~>
```

# Submitting Jobs

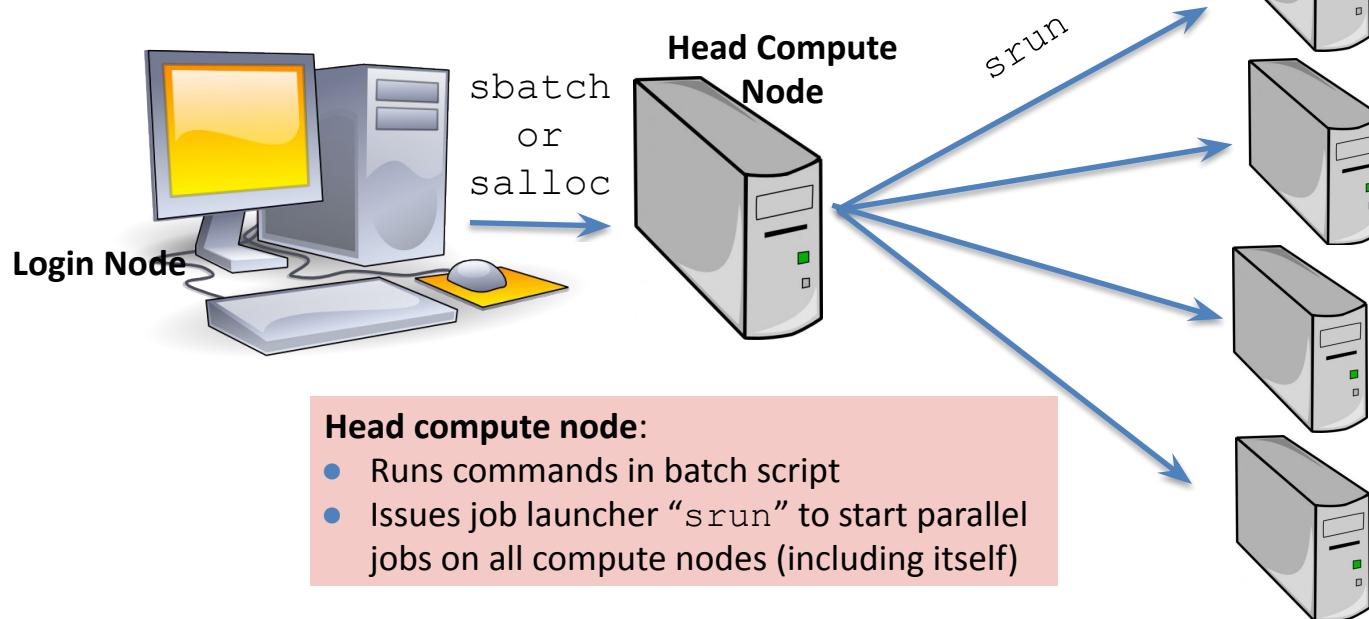
- Jobs can be submitted to queueing system through sbatch or salloc:
  - `sbatch <my_job_script>`
  - `salloc <options>`
- The above methods list details about resources needed for a job and for how long
- Eg.: a request for a cpu node to be allocated for 5 mins:
  - `salloc -N 1 -C cpu -q debug -t 5 -A <project>`



# Submitting Jobs

## Login node:

- Submit batch jobs via `sbatch` or `salloc`



\*figure courtesy Helen (2020 NERSC Training)

# My First “Hello World” Job Script

debug queue

2 nodes

10 min “walltime”

run on haswell partition

use SCRATCH file system

run 64 processes in parallel



my\_batch\_script.sh

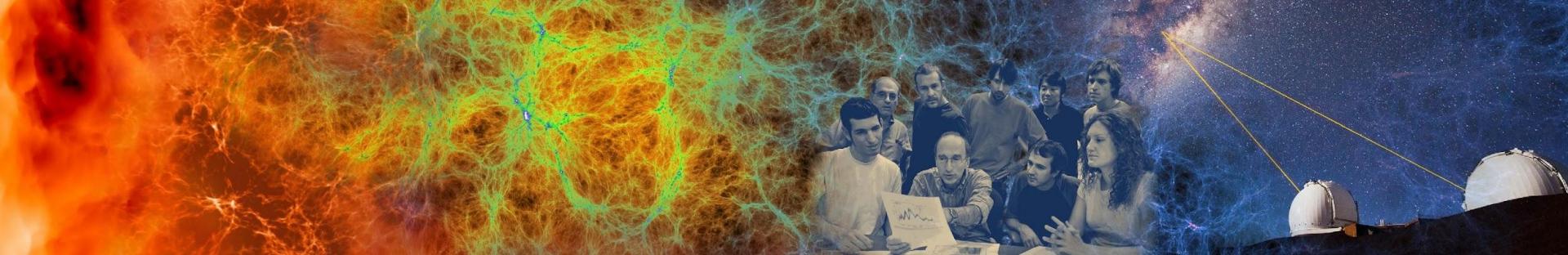
```
#!/bin/bash
#SBATCH -q debug
#SBATCH -N 2
#SBATCH -t 10:00
#SBATCH -C haswell
#SBATCH -L SCRATCH
#SBATCH -J myjob
srun -n 64 ./helloWorld
```

## To run via batch queue

```
% sbatch my_batch_script.sh
```

## To run via interactive batch

```
% salloc -N 2 -q interactive -C gpu -t 10:00
<wait_for_session_prompt. Land on a compute node>
% srun -n 64 ./helloWorld
```



# NERSC Training Accounts

for *The CUDA-Q Tutorial On Perlmutter*



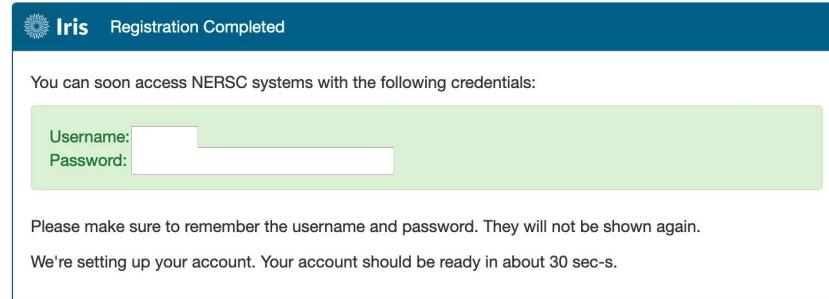
BERKELEY LAB



Office of  
Science

# Activate your Training Account

- Go to: <https://iris.nersc.gov/train>
- Enter the Training Code: “ep81” along with your details  
**Note:** if you have previously registered an account using a different code please recreate it.
- **Important: make a copy of your login and password**  
you won’t be able to change these, nor recover them later!!!



# Access to Perlmutter and Use Cuda-Q

- NERSC users have been added to ntrain10 project
- Non-users were sent the instruction to get a training account
  - Account valid through Nov 18th
- Login to Perlmutter: ssh [username@perlmutter.nersc.gov](mailto:username@perlmutter.nersc.gov)
- Cuda modules:
  - % module load cudatoolkit
- Running Jobs examples:
  - <https://docs.nersc.gov/jobs/>



# Compute Node Reservations

- GPU node reservation: Sun Nov 17th and Mon Nov 18th
  - To use 1 GPU only (sample flags for sbatch or salloc):
    - `-A ntrain1 --reservation=julia_sc24 -C gpu -N 1 -c 32 -G 1 -t 30:00 -q shared`
  - To use multiple nodes (sample flags for sbatch or salloc):
    - `-A ntrain1 --reservation=julia_sc24 -C gpu -N 2 -t 30:00 -q regular`
- Outside of reservation, use:
  - To use 1 GPU only (sample flags for sbatch or salloc):
    - `-A <project> -C gpu -N 1 -c 32 -G 1 -t 30:00 -q shared`
  - To use multiple nodes (sample flags for sbatch or salloc):
    - `-A <project> -C gpu -N 2 -t 30:00 -q regular (or -q interactive for salloc)`

# Reminder about file Systems

`$HOME` is not the best file system for running jobs at scale.  
Consider therefore replacing all occurrences of `$HOME` with  
`$SCRATCH`.

<https://github.com/NERSC/sc24-quantum-tutorial>

More detailed information in the github repository

# Logging into Jupyter

- Go to <https://jupyter.nersc.gov>
- Sign in using your training account credentials
- Select your preferred jupyter instance:

The screenshot shows the jupyterhub login page with the following details:

- Header:** jupyterhub Home Token Services Documentation train130 Logout
- MOTD (Maintenance Of The Day):** A yellow box contains the message: "The following NERSC resources that Jupyter depends on appear to be in maintenance or having issues. This may impact Jupyter. See the [NERSC MOTD](#) for further information." Below it, it says "Perlmutter status: degraded".
- Resource Selection Table:** A table comparing four node types:

	Shared CPU Node	Exclusive CPU Node	Exclusive GPU Node	Configurable GPU
<b>Perlmutter</b>	<a href="#">start</a>	<a href="#">start</a>	<a href="#">start</a>	<a href="#">start</a>
<b>Cori</b>	<a href="#">start</a>			<a href="#">start</a>
<b>Resources</b>	Use a node shared with other users' notebooks but outside the batch queues.	Use your own node within a job allocation using defaults.	Use multiple compute nodes with specialized settings.	
<b>Use Cases</b>	Visualization and analytics that are not memory intensive and can run on just a few cores.	Visualization, analytics, machine learning that is compute or memory intensive but can be done on a single node.	Multi-node analytics jobs, jobs in reservations, custom project charging, and more.	



# Setup

- Go to <https://jupyter.nersc.gov> and select:  
“Exclusive GPU Node” in the “Perlmutter” row

The screenshot shows the jupyterhub interface with a navigation bar at the top. Below the navigation bar is a yellow banner with a message about NERSC resources being in maintenance or having issues, and a status message for Perlmutter. The main content area displays a table comparing four resource types: Shared CPU Node, Exclusive CPU Node, Exclusive GPU Node, and Configurable GPU. The "Exclusive GPU Node" row for Perlmutter has a red box around its "start" button.

	Shared CPU Node	Exclusive CPU Node	Exclusive GPU Node	Configurable GPU
Perlmutter	<button>start</button>	<button>start</button>	<button>start</button>	<button>start</button>
Cori	<button>start</button>			<button>start</button>
Resources	Use a node shared with other users' notebooks but outside the batch queues.	Use your own node within a job allocation using defaults.		Use multiple compute nodes with specialized settings.
Use Cases	Visualization and analytics that are not memory intensive and can run on just a few cores.	Visualization, analytics, machine learning that is compute or memory intensive but can be done on a single node.		Multi-node analytics jobs, jobs in reservations, custom project charging, and more.



# Setup

If you're spending a long time in the queue waiting for an exclusive GPU node, "Shared GPU Node"s are also available – eg:

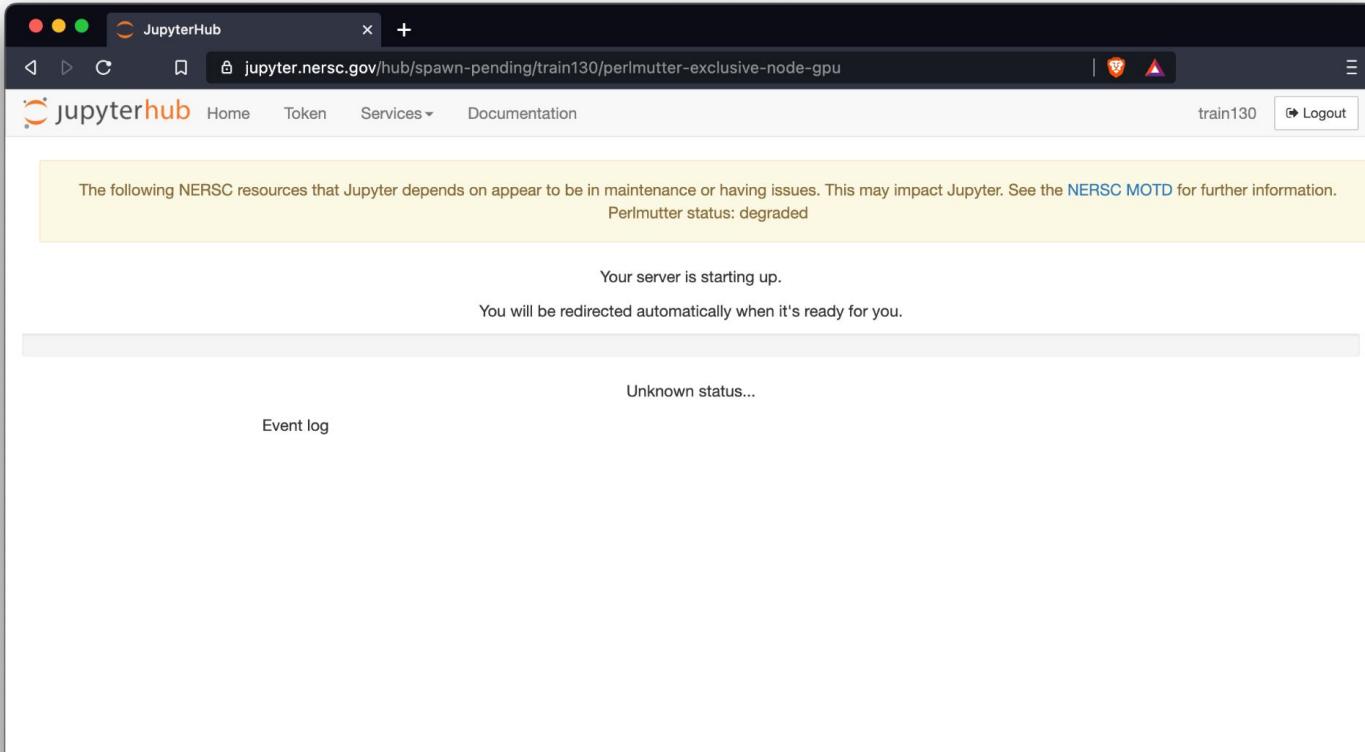
	Login Node	Shared GPU Node	Exclusive CPU Node	Exclusive GPU Node	Configurable Job
Alvarez	<button>start</button>	<button>start</button>	<button>start</button>	<button>start</button>	
Muller	<button>start</button>	<button>start</button>	<button>start</button>	<button>start</button>	<button>start</button>
Perlmutter	<button>stop</button> <button>server</button>	<button>start</button>	<button>start</button>	<button>start</button>	<button>start</button>
Resources	Use a login node shared with other users, outside the batch queues.	Use a single GPU on a node within a job allocation using defaults.	Use your own node within a job allocation using defaults.	Use multiple compute nodes with specialized settings.	
Use Cases	Visualization and analytics that are not memory intensive and can run on just a few cores.	Work that fits on a single GPU, and uses at most a quarter of a GPU node's CPU cores and host memory.	Visualization, analytics, machine learning that is compute or memory intensive but can be done on a single node.	Multi-node analytics jobs, jobs in reservations, custom project charging, and more.	

Alternatively please use a reservation (details later)



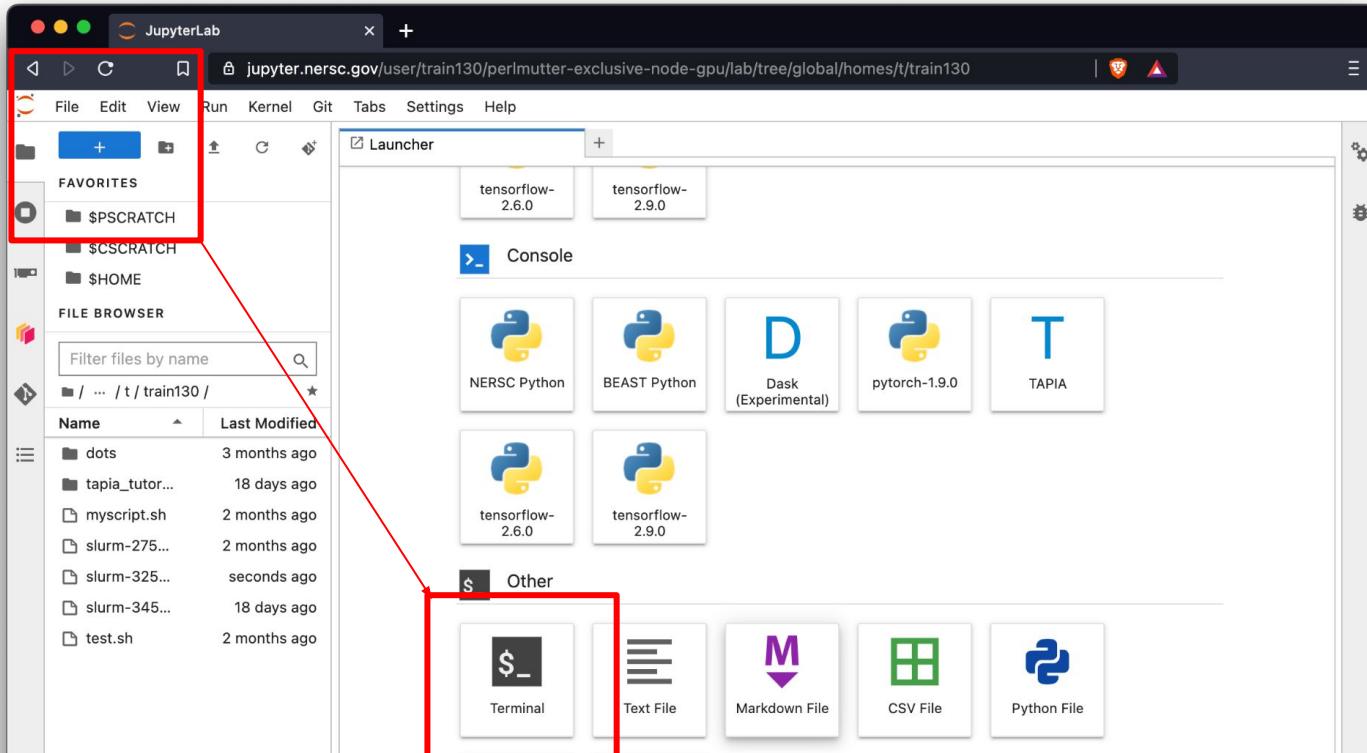
# Setup

- Jupyter should take a minute to start:



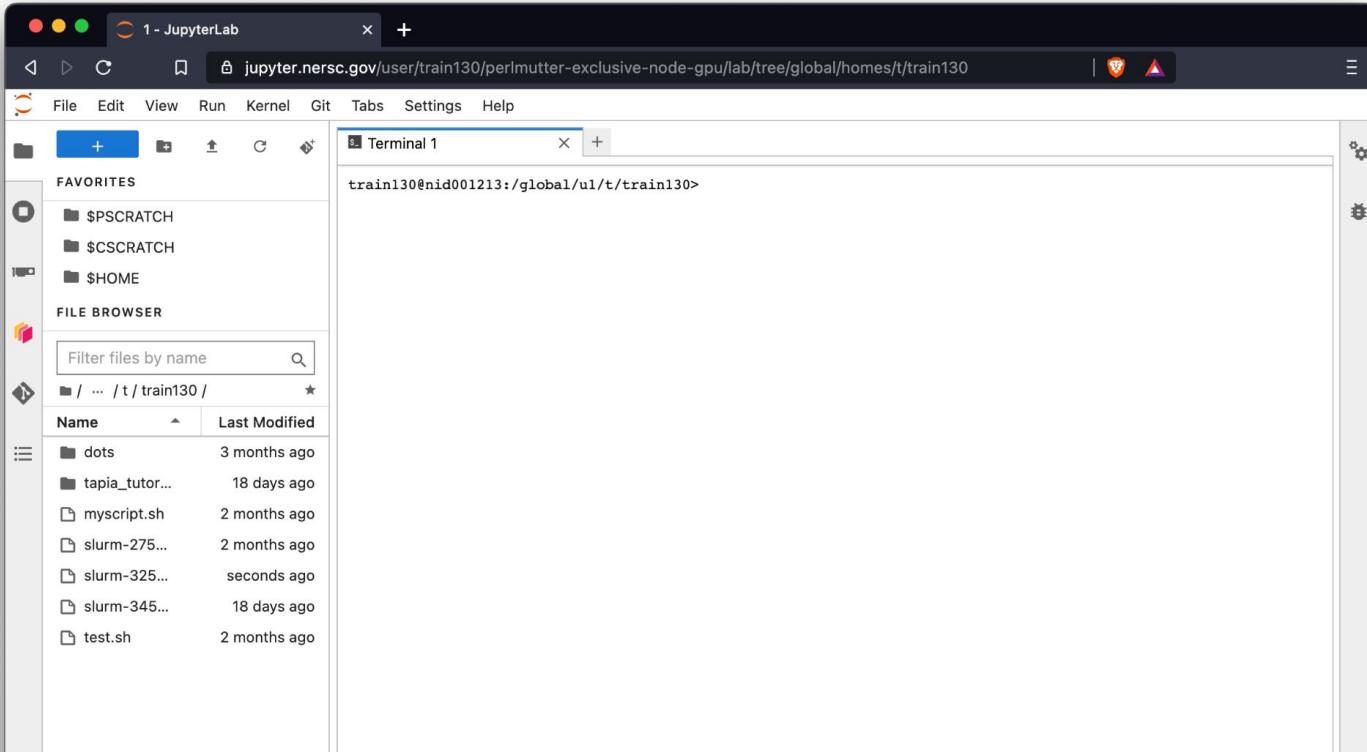
# Setup

- (If you don't see a terminal), select “+” followed by “Terminal”



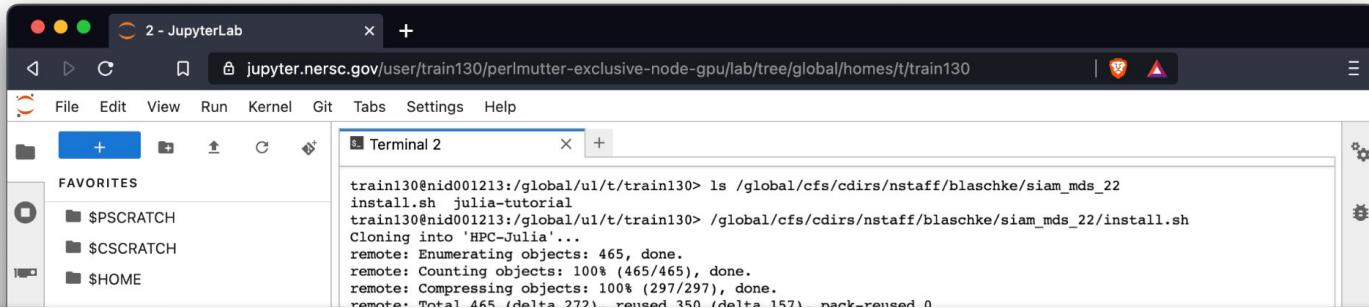
# Setup

- Once started you should see a terminal



# Setup

- (In the terminal), follow installation instructions



A screenshot of a JupyterLab interface on a Mac OS X system. The title bar says "2 - JupyterLab". The left sidebar shows "FAVORITES" with entries for "\$PSCRATCH", "\$CSCRATCH", and "\$HOME". The main area has a terminal tab titled "Terminal 2". The terminal window displays the following command-line session:

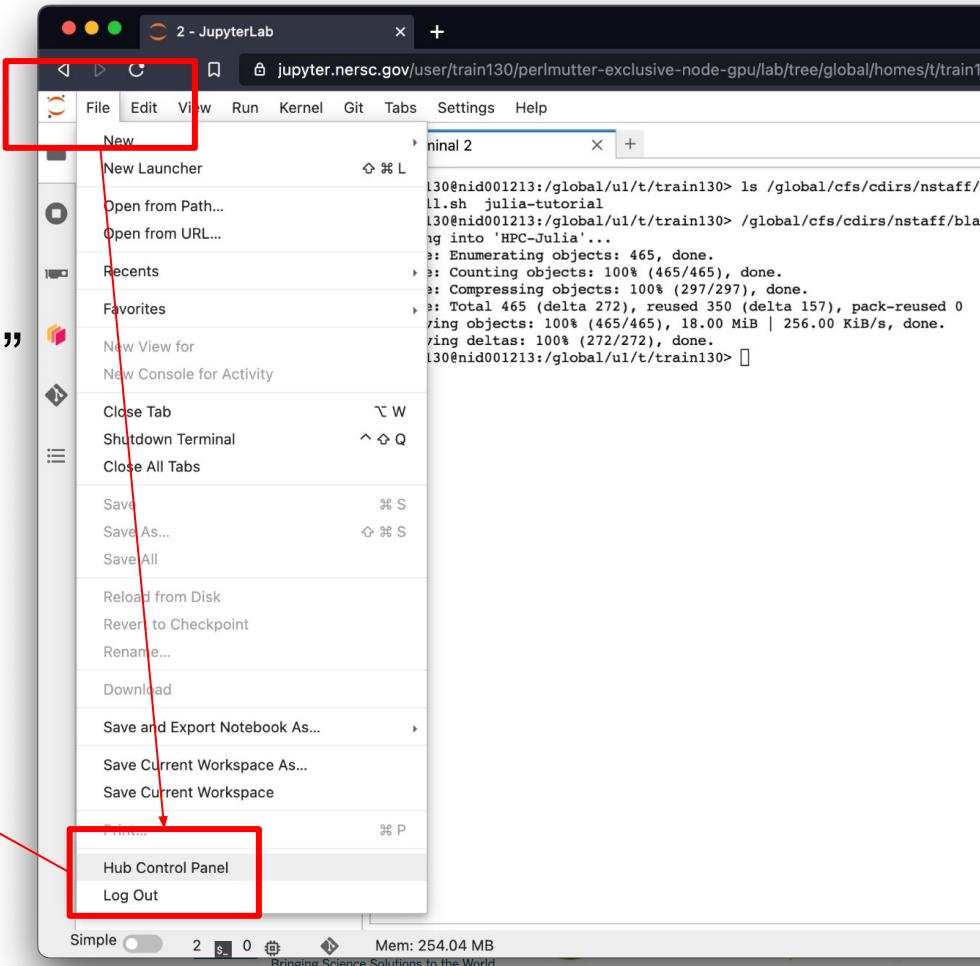
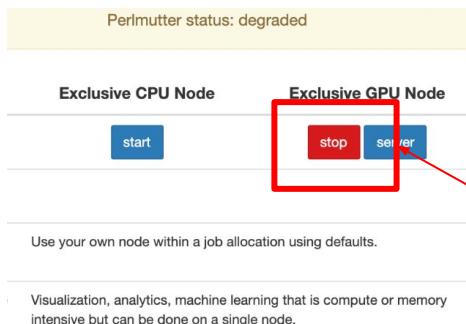
```
train130@nid001213:/global/u1/t/train130> ls /global/cfs/cdirs/nstaff/blaschke/siam_mds_22
install.sh julia-tutorial
train130@nid001213:/global/u1/t/train130> /global/cfs/cdirs/nstaff/blaschke/siam_mds_22/install.sh
Cloning into 'HPC-Julia'...
remote: Enumerating objects: 465, done.
remote: Counting objects: 100% (465/465), done.
remote: Compressing objects: 100% (297/297), done.
remote: Total 465 (delta 272), reused 350 (delta 157), pack-reused 0
```

Upload the folders and script to your scratch directory:

```
$ scp folder_name username@perlmutter.nersc.gov:/pscratch/your_path
```

# Setup

- Stop your server by going to:  
“File” > “Hub Control Panel”
- The selecting “Stop”



# Setup

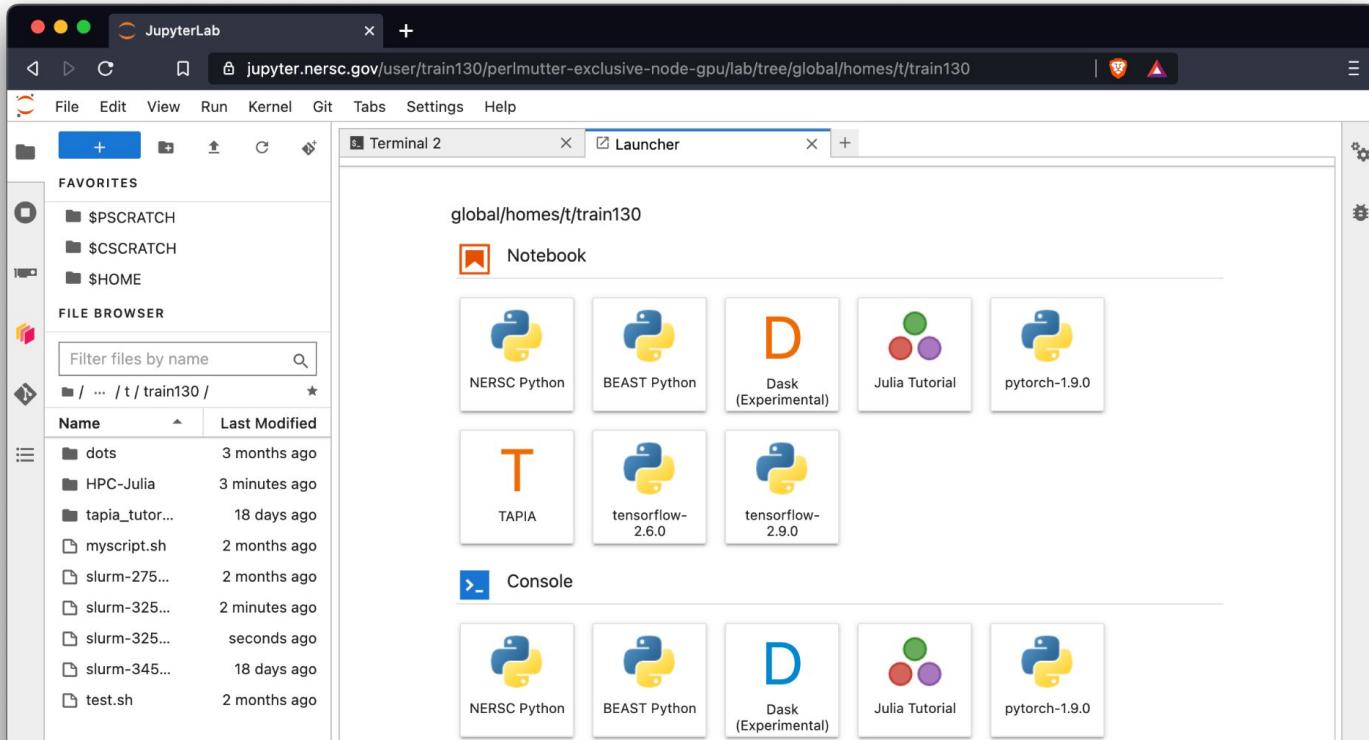
- Start a fresh exclusive node on Perlmutter

The following NERSC resources that Jupyter depends on appear to be in maintenance or having issues. This may impact Jupyter. See the [NERSC MOTD](#) for further information.  
Perlmutter status: degraded

	Shared CPU Node	Exclusive CPU Node	Exclusive GPU Node	Configurable GPU
<b>Perlmutter</b>	<a href="#">start</a>	<a href="#">start</a>	<a href="#">start</a>	<a href="#">start</a>
<b>Cori</b>	<a href="#">start</a>			<a href="#">start</a>
<i>Resources</i>	Use a node shared with other users' notebooks but outside the batch queues.	Use your own node within a job allocation using defaults.	Use multiple compute nodes with specialized settings.	
<i>Use Cases</i>	Visualization and analytics that are not memory intensive and can run on just a few cores.	Visualization, analytics, machine learning that is compute or memory intensive but can be done on a single node.	Multi-node analytics jobs, jobs in reservations, custom project charging, and more.	

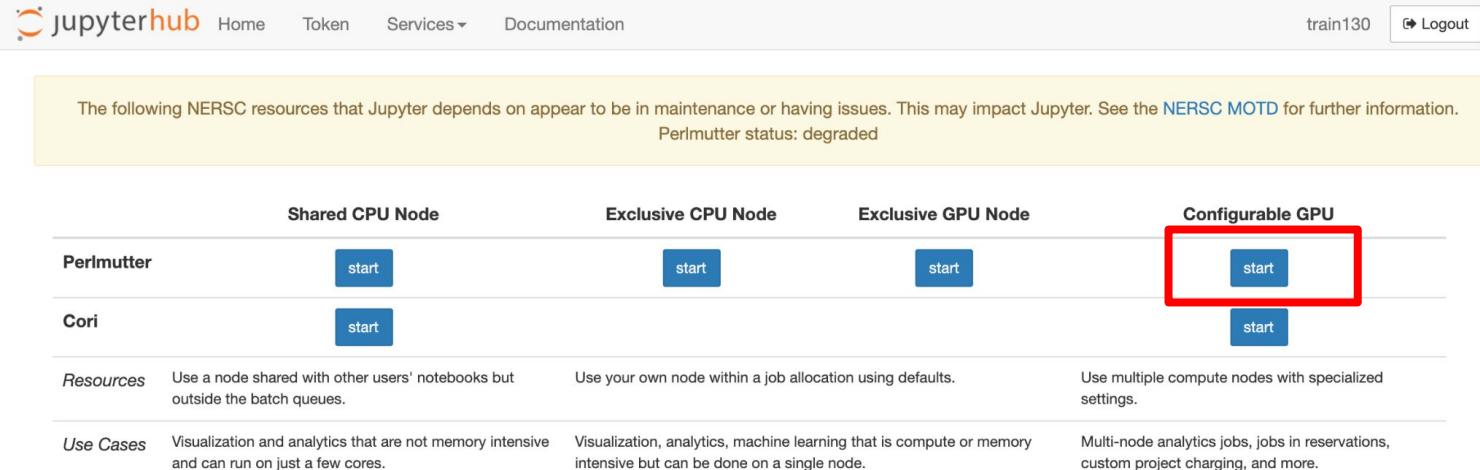
# Setup

- Select the launcher, you should see the Cuda-Q Kernel :)



# Using Reservations in Tutorials

- Go to <https://jupyter.nersc.gov> and select: “Configurable GPU” in the “Perlmutter” row



The following NERSC resources that Jupyter depends on appear to be in maintenance or having issues. This may impact Jupyter. See the [NERSC MOTD](#) for further information.  
Perlmutter status: degraded

	Shared CPU Node	Exclusive CPU Node	Exclusive GPU Node	Configurable GPU
<b>Perlmutter</b>	<a href="#">start</a>	<a href="#">start</a>	<a href="#">start</a>	<a href="#">start</a>
<b>Cori</b>	<a href="#">start</a>			<a href="#">start</a>
<b>Resources</b>	Use a node shared with other users' notebooks but outside the batch queues.	Use your own node within a job allocation using defaults.		Use multiple compute nodes with specialized settings.
<b>Use Cases</b>	Visualization and analytics that are not memory intensive and can run on just a few cores.	Visualization, analytics, machine learning that is compute or memory intensive but can be done on a single node.		Multi-node analytics jobs, jobs in reservations, custom project charging, and more.



# Jupyter Options:

Leave defaults, except:

Account:

ntrain10

Reservation

TBD

Time

TBD

The screenshot shows a web browser window titled "JupyterHub" with the URL "jupyter.nersc.gov/hub/spawn/train130/perlmutter-configurable-gpu". The page displays a message about NERSC resources being in maintenance or having issues, specifically mentioning "Perlmutter status: degraded". Below this, the title "Server Options" is centered. The configuration form includes the following fields:

- Account ("\_g" suffix will be added as needed): ntrain8
- Constraint: gpu
- QOS: jupyter
- cpus-per-task (node has 128 cpus): 128
- gpus-per-task (node has 4 GPUs): 4
- nodes (maximum of 4 for jupyter QOS): 1
- ntasks-per-node: 1
- Reservation: siam\_intro\_gnn
- time (time limit in minutes): 120

A large orange "Start" button is located at the bottom of the form.



# Note:

- The reservation is active from until 5pm EST today. If you request a node during this time, make sure that the current time + requested time does not exceed the reservation's end time (5pm).
- Training accounts are fully fledged NERSC accounts
- Will remain available until Tuesday – please copy any data and code to your computer before then



# Agenda

## Morning Sessions

- Introductions
- Overview of CUDA-Q
- Onboarding to Perlmutter
- Accelerating Quantum Simulations Notebook
- NVIDIA Quantum Cloud and Local install of CUDA-Q

## Afternoon Sessions

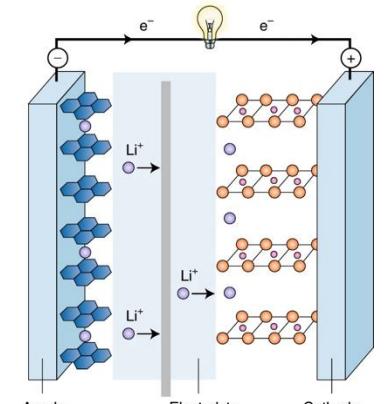
- Intro to large scale clusters, how to navigate and use them
- Example: Quantum chemistry at NERSC
- Industry Use Case: Simulating Hamiltonians of molecules with 30,000 terms
- QuEra's Quantum Hardware
- Example: Quantum reservoir computing

# What's the problem?

- The goal is finding the ground state of a many-body Hamiltonian (describing e.g., a many-electron molecule)

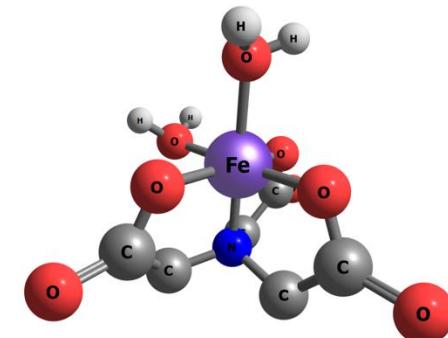
$$H = \sum_i \left( -\frac{\hbar^2}{2m} \nabla_i^2 + V_{\text{ext}}(r_i) \right) + \sum_{i,j} V_{\text{e-e}}(|r_i - r_j|)$$

External potential (nuclei-electrons)  
Kinetic part for the electrons  
Interaction (electrons-electrons)



Nat. Electron.

- In general,  $H = H_1 + H_2$ 
  - $H_1$  contains only one-body operators (kinetic terms, external potentials, ...). **Easy**
  - $H_2$  contains the two-body operators (Coulomb interactions, ...). **Complicated**
- Solving for the ground state is a difficult task:
  - Space of solutions grows exponentially
  - Perturbation theory often fails
- Need for numerical / approximate techniques



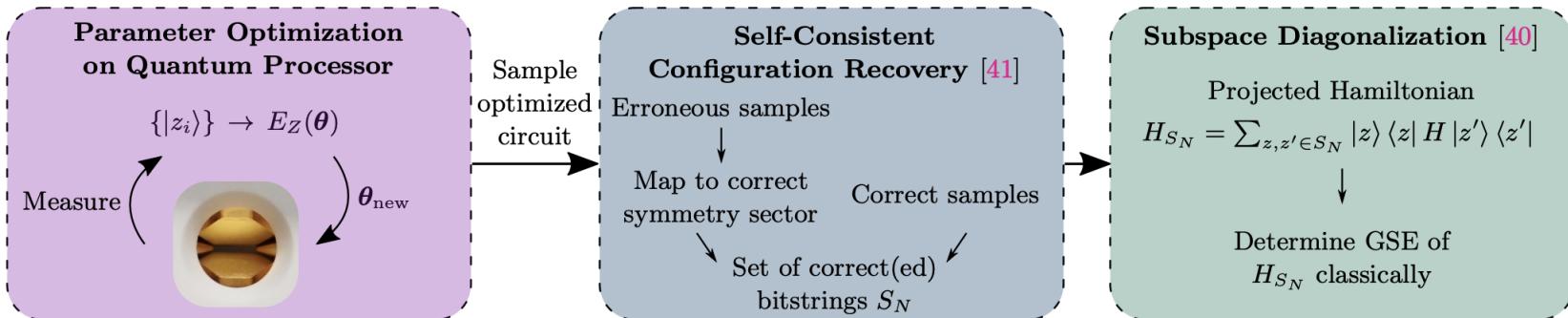
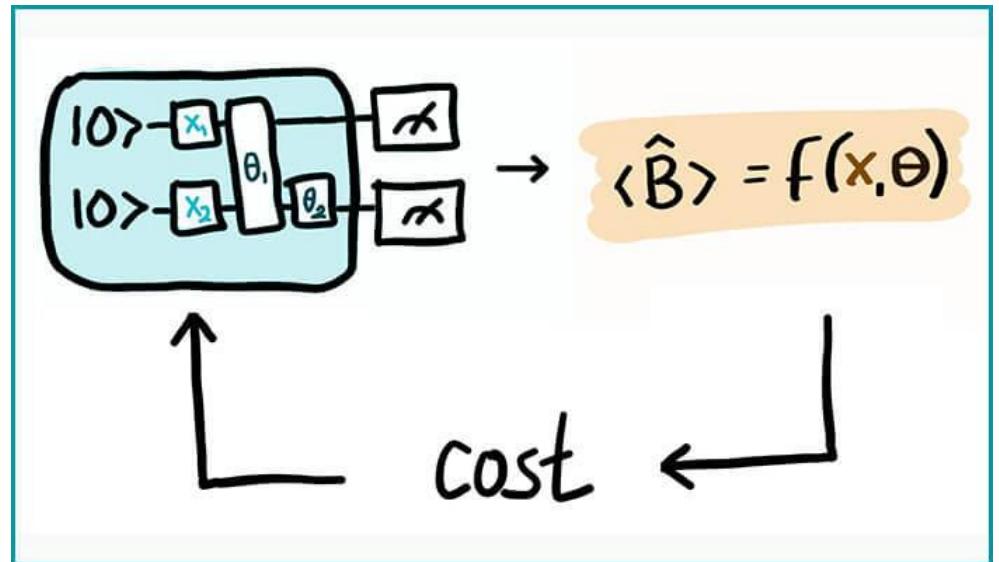
arxiv:2408.1080

1

# VQE: One possible solution

- In NISQ devices, VQE with hybrid approaches:

$$E = \min_{\theta} \langle \psi(\theta) | H | \psi(\theta) \rangle$$



Solving an Industrially Relevant Quantum Chemistry Problem on Quantum Hardware

Ludwig Nützel,<sup>1,\*</sup> Alexander Gresch,<sup>2,3</sup> Lukas Hehn,<sup>4</sup> Lucas Marti,<sup>1</sup> Robert Freund,<sup>5</sup> Alex Steiner,<sup>5</sup> Christian D. Marciak,<sup>5</sup> Timo Eckstein,<sup>1,6</sup> Nina Stockinger,<sup>1,7</sup> Stefan Wolf,<sup>1</sup> Thomas Monz,<sup>5,†</sup> Michael Kühn,<sup>4</sup> and Michael J. Hartmann<sup>1,6</sup>

<sup>1</sup>Department of Physics, Friedrich-Alexander Universität Erlangen-Nürnberg, Erlangen, Germany

<sup>2</sup>Heinrich Heine University Düsseldorf, Faculty of Mathematics and Natural Sciences, Germany

<sup>3</sup>Institute for Quantum Inspired and Quantum Optimization, Hamburg University of Technology, Germany

<sup>4</sup>BASF SE, Next Generation Computing, Pfalzgrafenstr. 1, 67061 Ludwigshafen, Germany

<sup>5</sup>Universität Innsbruck, Institut für Experimentalphysik, Innsbruck, Austria

<sup>6</sup>Max Planck Institute for the Science of Light, Staudtstraße 2, 91058 Erlangen, Germany

<sup>7</sup>Department of Chemistry, Technical University of Munich, Garching, Germany

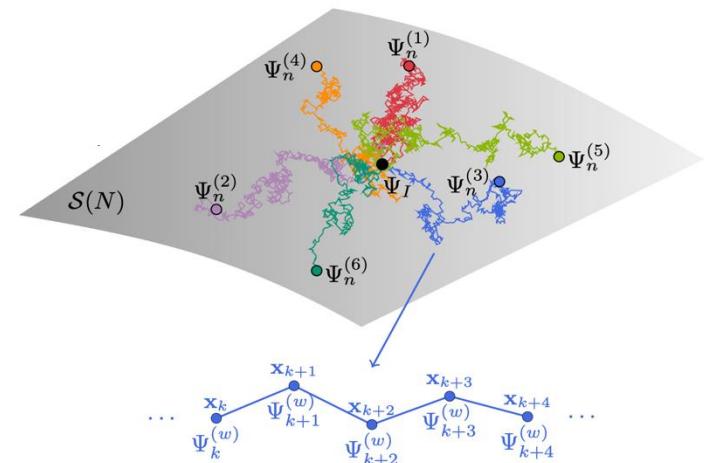
# Auxiliary Field QMC

M Motta, S Zhang, WIREs Comput Mol Sci 2018, 8:e1364

- AFQM transforms the “complicated” interaction terms to an infinite sets of non-interacting terms
- Strongly correlated systems not well described by lower level of theory (HF, DFT, ...)
- Monte Carlo methods proved useful in many other applications

## AFQMC: basic ingredients

- Imaginary time evolution
- Random walk in the space of Slater determinants
- Good trial (initial) wave function



# AFMQC: Imaginary time evolution

AFQMC is a projection quantum Monte Carlo technique

- Our target is the ground state  $|\Psi_0\rangle$  of  $H$

- Let's consider an initial state

$$|\Phi_I\rangle = c_0|\Psi_0\rangle + c_1|\Psi_1\rangle + c_2|\Psi_2\rangle + \dots$$

Excited states of  $H$

- Apply  $e^{-\beta H}$

$$E_0 < E_1 < E_2 < \dots$$

Eigenvalues of  $H$

$$e^{-\beta H}|\Phi_I\rangle \sim c_0|\Psi_0\rangle + c_1 e^{-\beta(E_1-E_0)}|\Psi_1\rangle + c_2 e^{-\beta(E_2-E_0)}|\Psi_2\rangle + \dots$$

- All the high-energy states are exponentially suppressed
- We project onto the ground state when  $\beta \rightarrow \infty$ .

# Why do we need a quantum computer?

- For extracting physical quantities like the g.s. energy, a **trial wf**  $|\Psi_T\rangle$  is generally used

$$E_0 = \langle \Psi_T | H e^{-\beta H} | \Phi_I \rangle$$

- $|\Psi_T\rangle$  can be a high quality wf hard to simulate on a classical computer, e.g.:

$$|\Psi_T\rangle = |\Psi_{CC}\rangle = \exp(T - T^\dagger) |\phi_0\rangle$$

(unitary coupled cluster)

$$T = \sum_{P,Q} t_{PQ} \hat{a}_P^\dagger \hat{a}_Q + \sum_{P,Q,R,S} t_{PQ}^{RS} \hat{a}_P^\dagger a_Q^\dagger \hat{a}_R \hat{a}_S$$

- What do we need from the QC?

$$c = \langle \Psi_T | e^{-\beta H} | \Phi_I \rangle$$

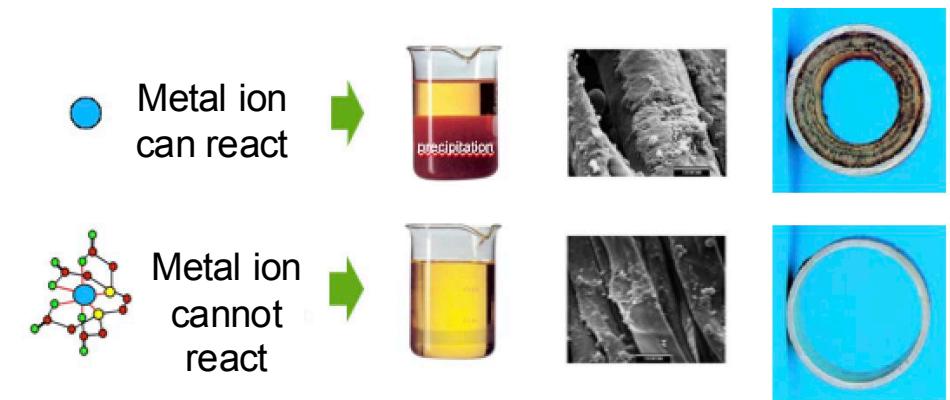
overlaps

Computed via Hadamard test or shadow techniques

# FeNTA chelating agent for industrial applications

- Chelating agents are used to shield the undesired influence of metal ions in various industrial processes and environmental systems.

- They are used
  - in the detergent industry, to remove hard water ions and improve cleaning ability of detergents;
  - in paper and pulp production, to improve the effectiveness of the bleaching chemicals;
  - in the mining industry, they form stable complexes with metal ions, making it easier to separate and extract valuable metals such as copper, gold, and silver from the ore;
  - in the textile industry to remove metal ions from fabrics, preventing discoloration and damage, and improving the dyeing process;
  - in the photographic industry, to avoid unpleasant deposits on the photographic layers and contamination of the developing equipment;
  - in the agricultural industry, they continuously supply the plants with very low concentrations of metal ions required for healthy growth.



BASF has been producing these chemicals for many decades. The range and application possibilities grew continuously in line with customer requirements. Today, BASF is one of the leading international producers of organic chelating agents.

# Agenda

## Morning Sessions

- Introductions
- Overview of CUDA-Q
- Onboarding to Perlmutter
- Accelerating Quantum Simulations Notebook
- NVIDIA Quantum Cloud and Local install of CUDA-Q

## Afternoon Sessions

- Intro to large scale clusters, how to navigate and use them
- Example: Quantum chemistry at NERSC
- Industry Use Case: Simulating Hamiltonians of molecules with 30,000 terms
- QuEra's Quantum Hardware
- Example: Quantum reservoir computing



# Applications of Neutral-Atom Analog Quantum Computing

Tommaso Macrì, QuEra Computing



November 2024

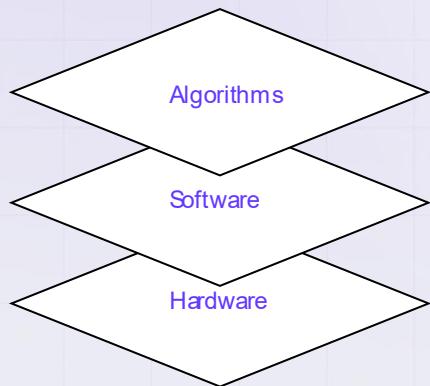
# About QuEra

- Based in Boston, we employ 50+ scientists and engineers from Harvard, MIT, and other top-tier institutions.
- We have a working analog quantum computer on a public cloud and our mission is to build and commercialize **large-scale, error-corrected digital devices**.
- We productize innovations from QuEra and our MIT/Harvard collaborators at **record speed**.
- Dozens of organizations and hundreds of users already execute jobs on our quantum computer.
- We are hiring, check current job opportunities!



# QuEra: A Full-Stack Approach

Tightly integrated hardware, software, and applications



To leverage unique hardware features and maximize the quantum compute power



Unique algorithms for optimization, machine learning, and simulation



Development and emulation environment to test quantum software: Available in Julia and Python

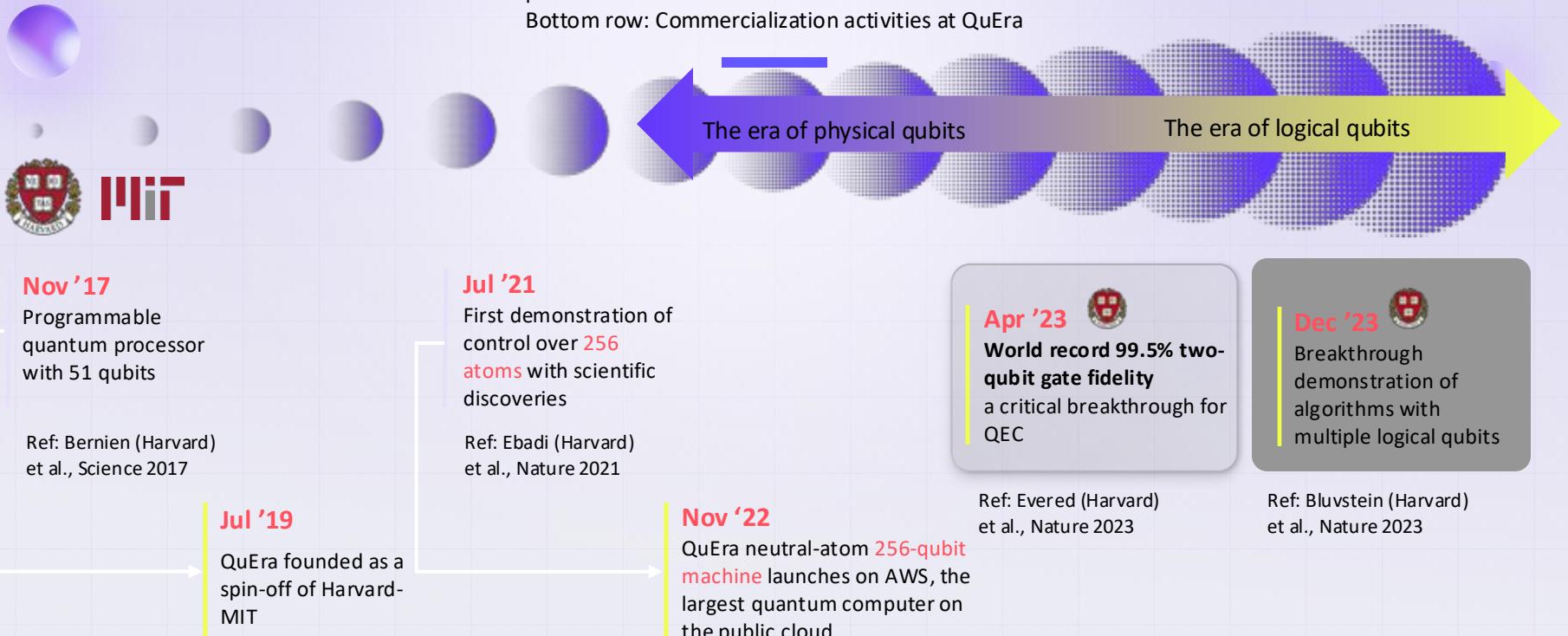


Flexible, large-scale quantum hardware with unique features

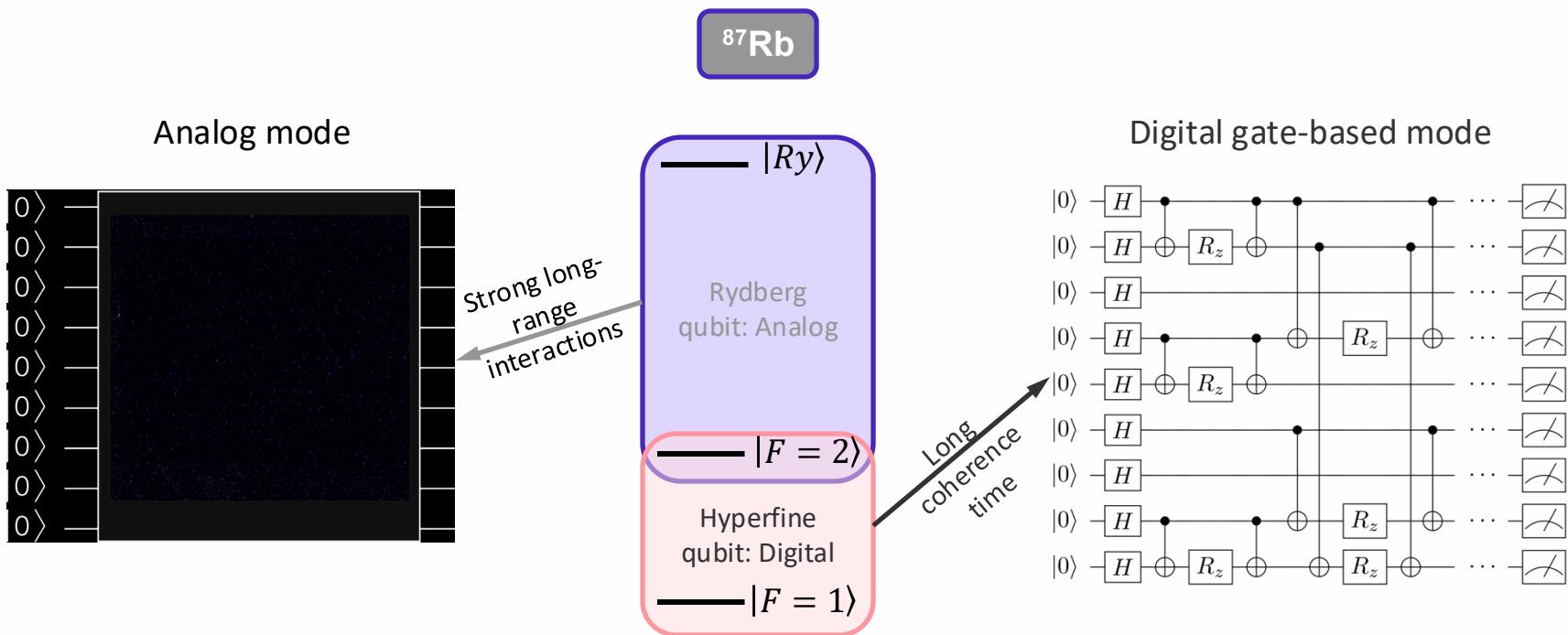
# Transitioning From Physical to Logical Qubits

Top row: Academic demonstrations at Harvard and MIT

Bottom row: Commercialization activities at QuEra



# Analog and Digital Operation



# Analog quantum computing system

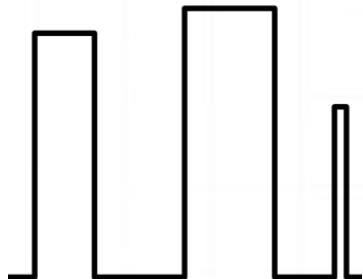
## Analog waveforms for time-dependent dynamics

User-defined waveforms execute many different protocols.

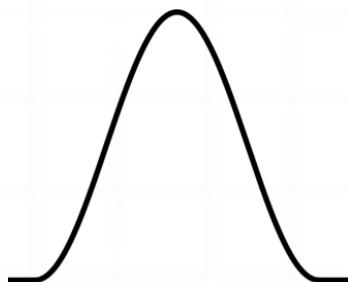
$$\hat{H}/\hbar = \sum_i \frac{\Omega(t)}{2} \left( e^{i\varphi(t)} |1\rangle_i \langle r| + h.c. \right) - \Delta_i(t) n_i + \sum_{i < j} V_{ij} n_i n_j$$

## Programmable connectivity

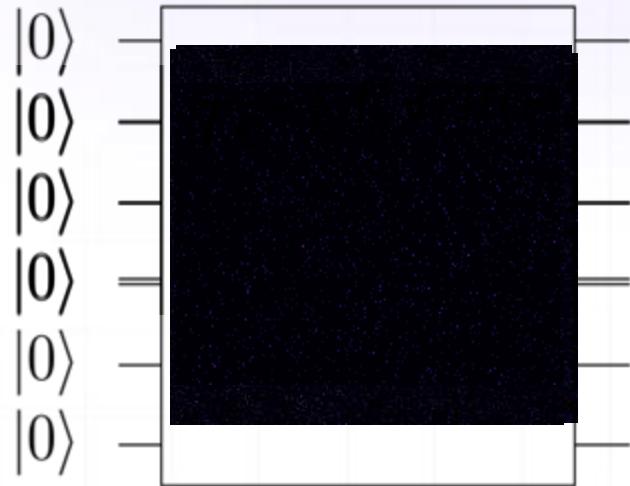
Encode problems in the spatial positions of atoms



Gates



Annealing



Arbitrary

# Aquila: Our 256-Qubit Analog Quantum Computer

- The first and still the world's **only** publicly-accessible quantum computer based on neutral atoms.
- Increased availability from 10 to over **100 hours per week**.
- Key applications: **simulation, optimization, machine learning**.
- The use of Aquila has resulted in several exciting scientific publications.
- Through "premium access", customers can engage directly with QuEra scientists.



# Overview of New Features & Outlooks

New Aquila features for Analog Hamiltonian Simulation (AHS)

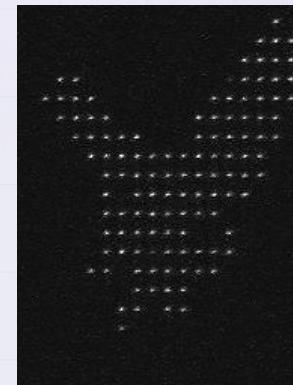
- Local detuning
  - Programmable individual qubit detunings in AHS
- More flexible qubit arrangements
  - Tighter row spacing 4->2 um, larger height 75 ->100 um

New opportunities in **digital quantum simulation**:

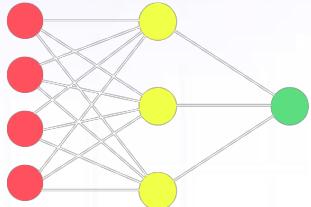
- Material, chemistry, condensed matter system
- Qubit shuttling, long-range connectivity
- Large scale
- Early fault tolerant digital quantum simulation
- Quantum advantage?
- ...

New opportunities in analog **quantum simulation**:

- Programmable disorder, initial state
- Non-equilibrium dynamics
- Exotic lattices
- Larger system sizes for 1D and quasi-1D problems



# What problems can QuEra's Aquila solve?



## Machine learning

- Classification
- Regression and prediction
- Reservoir computing



## Optimization

- Maximum independent set
- Non-native optimization, e.g. QUBO
- Commercially-relevant problems



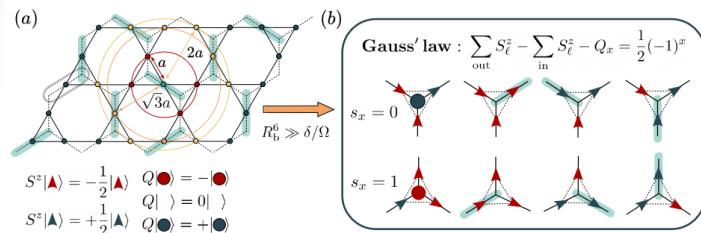
## Simulation

- Equilibrium quantum phase of matter
- Non-equilibrium quantum dynamics
- Lattice gauge theories

Check out the Aquila Whitepaper: Wurtz et al., arXiv:2306.11727

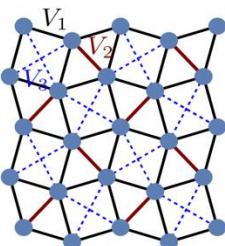
# Quantum Simulations Demonstrations on Aquila

## String Breaking



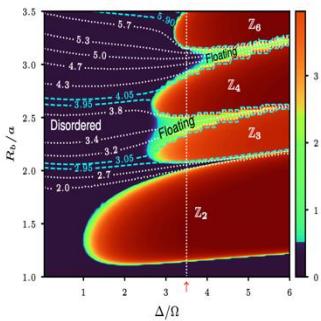
Collab. with Peter Zoller (Innsbruck)

## Shastry-Sutherland



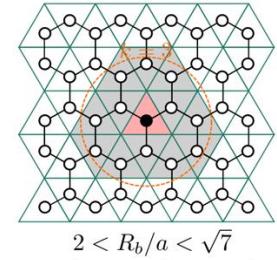
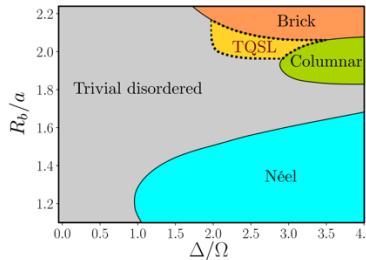
Collab. With Shalaev, Boltasseva, Banerjee groups  
(Purdue & QSC)

## Quantum Floating Phases



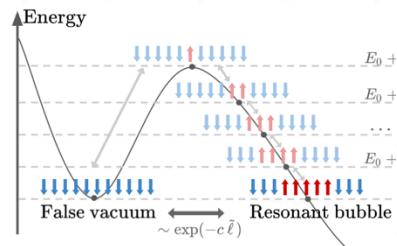
Collab. with Yannick Meurice  
(Iowa), Shan-Wen Tsai (UC  
Riverside), arXiv:2401.08087

## Trimer Quantum Spin Liquids



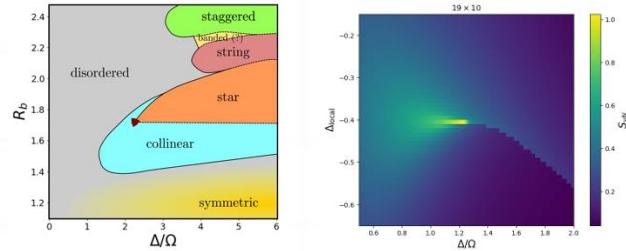
M. Kornjaca et al., Comm. Phys. 2023

## False Vacuum Dynamics



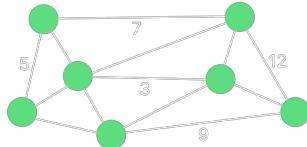
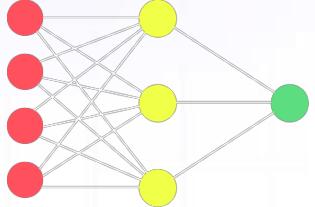
Collab. with NERSC  
S. Darbha et al, PRB (2024)

## Multi-Criticality and Analog Higgs



Collab. with NERSC  
S. Darbha et al, PRB (2024)

# What problems can QuEra's Aquila solve?



## Machine learning

- Classification
- Regression and prediction
- Reservoir computing

## Optimization

- Maximum independent set
- Non-native optimization, e.g. QUBO
- Commercially-relevant problems

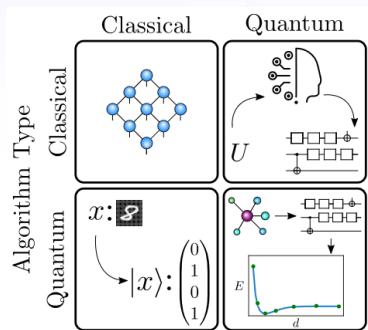
## Simulation

- Equilibrium quantum phase of matter
- Non-equilibrium quantum dynamics
- Lattice gauge theories

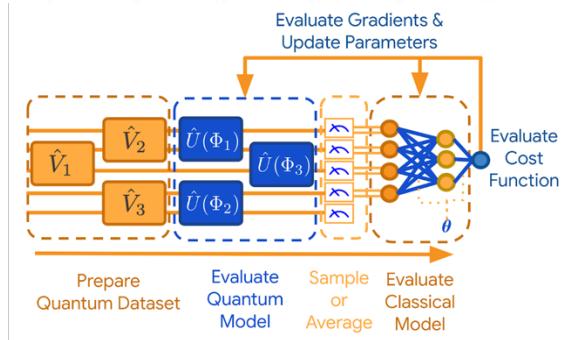
Check out the Aquila Whitepaper: Wurtz et al., arXiv:2306.11727

# Overview

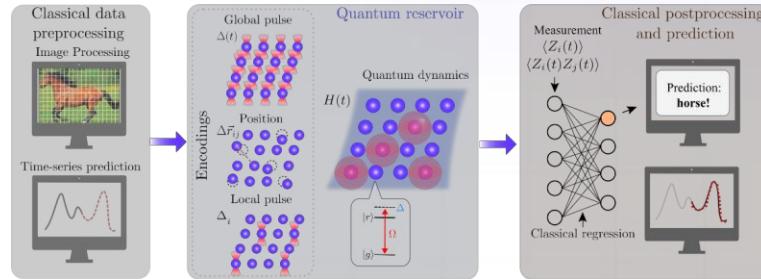
- Why quantum machine learning?



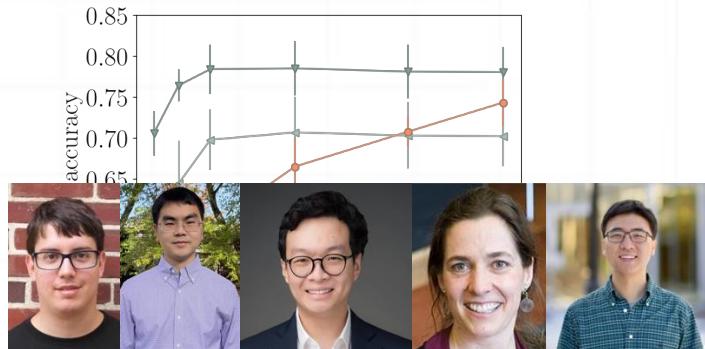
- Current quantum algorithms.



- Our approach – quantum reservoir learning.

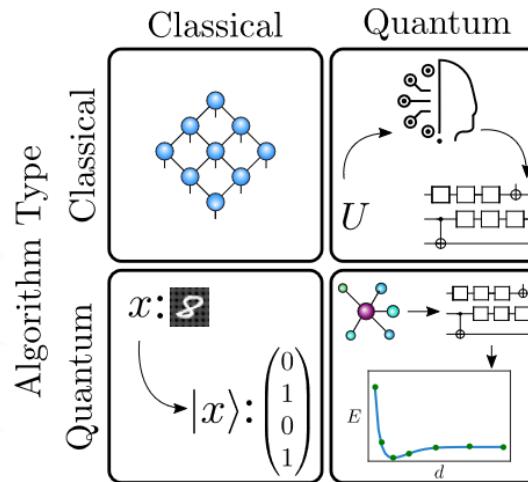


- Results – largest quantum machine learning experiments (108 qubits), kernel advantage.



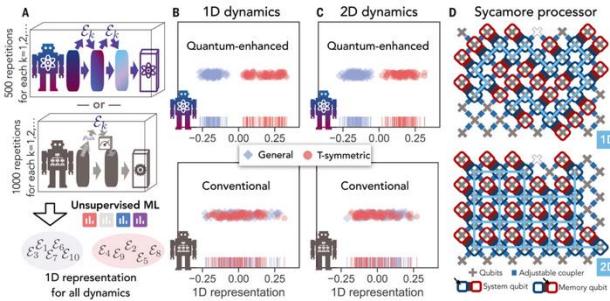
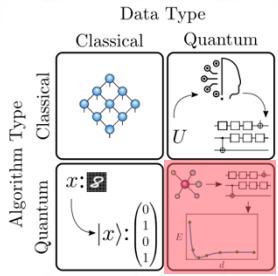
# Quantum machine learning

- Machine learning – “unreasonably” successful **data-based** computation paradigm.
- Quantum machine learning – **critical** to the future of the field.
- Prospect for advantage – **quantum data**.



# Quantum machine learning – promise

- Explicit – learning quantum dynamics.

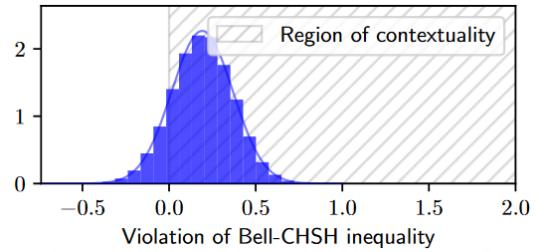


[Science 376, 6598 (2022)] \*40 qubits

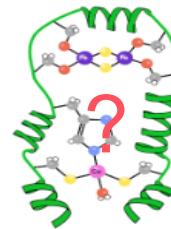
- Potential – materials and drug properties?

- Implicit – language ambiguities.

*A and B belong to the same [cannibalistic / herbivorous]<sub>1</sub> species of animal. In a hot afternoon in south Sahara, one of them<sub>1</sub> was very hungry. They notice each other when they were roaming in the field. In a while, one of them<sub>2</sub> is no longer [hungry / alive]<sub>2</sub>.*



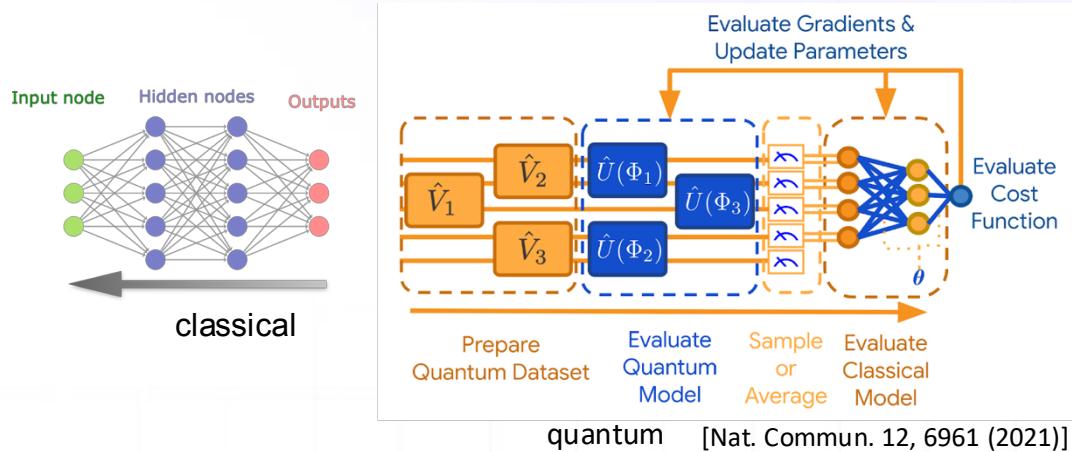
[EPTCS 384, 187 (2023)]



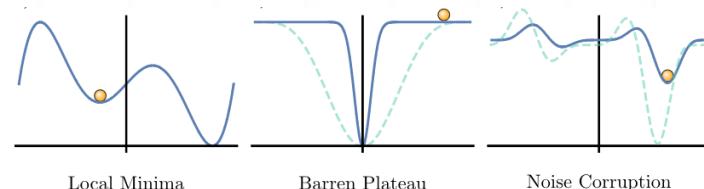
- What about **algorithms**?

# Quantum machine learning – practice

- **Variational QML** - quantum neural network.
- Dominate theoretical/numerical explorations.



- Major issues:
  - **Cost** – gradient estimation loop.
  - **Trainability** – barren plateaus.

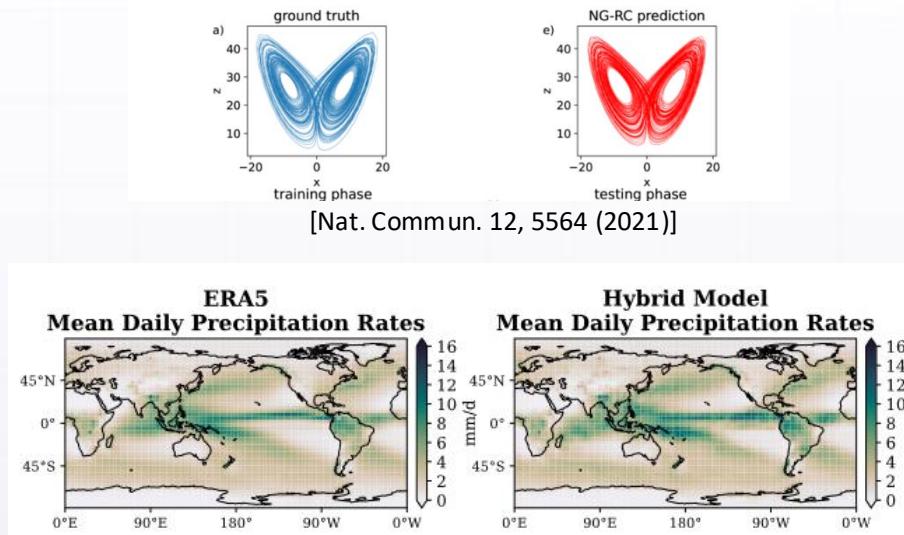
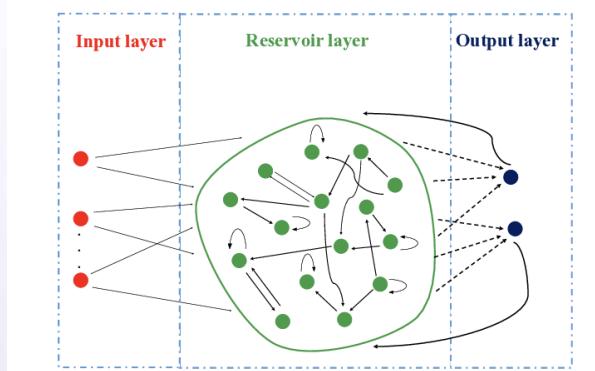


[Nat. Comput. Sci. 2, 567 (2022)]

- **Experiments** – rare, small-scale, limited advantage potential.
- Our approach – **quantum reservoir computing**.

# Classical reservoir computing

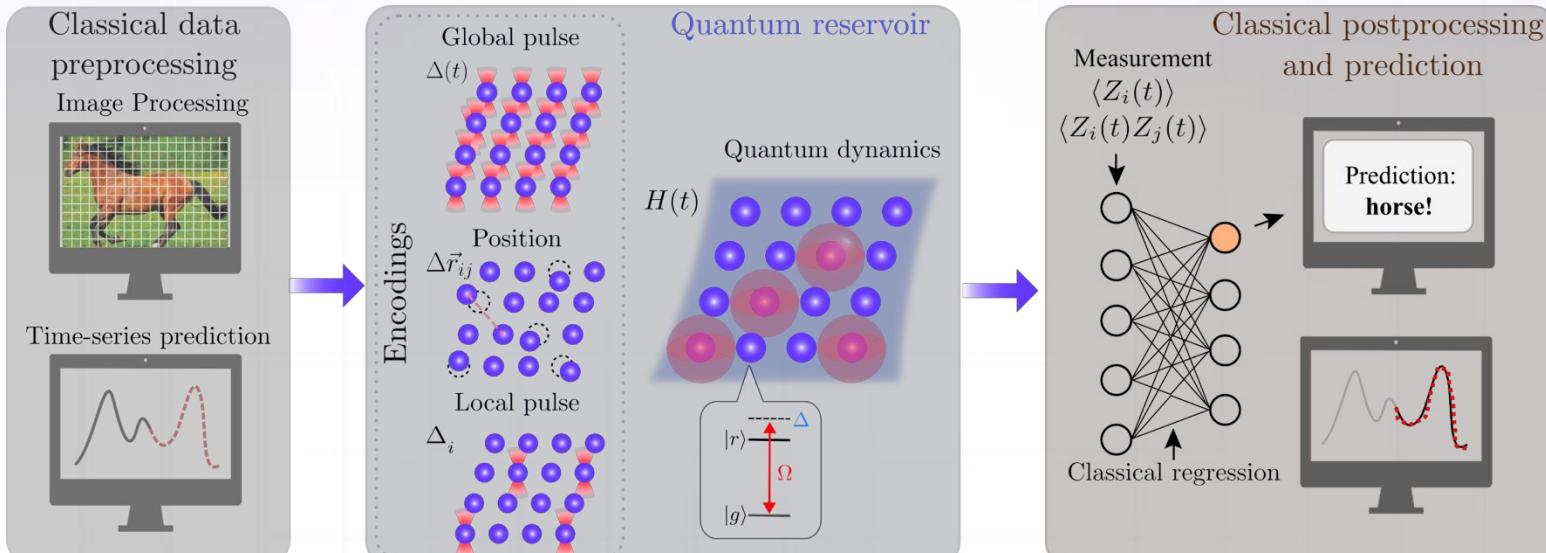
- Predictive power from a dynamical system – a **reservoir**.
- **Gradient-free** training.
- **Successful applications:**
  - Chaotic dynamics prediction.
  - Hybrid deployment – climate modeling.



[Geophys. Res. Lett. 50, 102649 (2023)]

# Quantum reservoir computing

Neutral atom quantum reservoir architecture:



## QRC advantages (over variational QML):

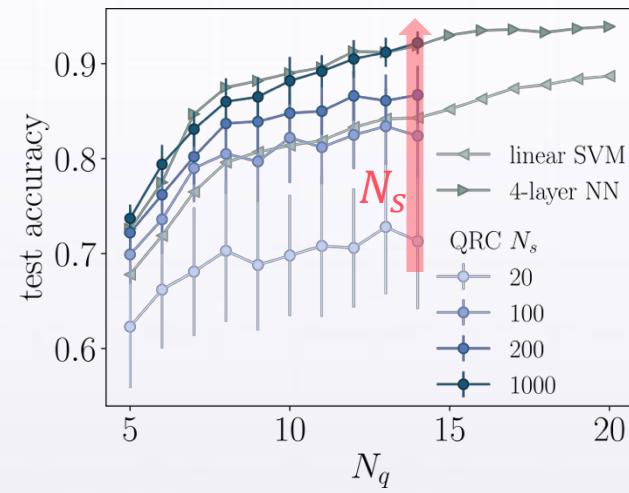
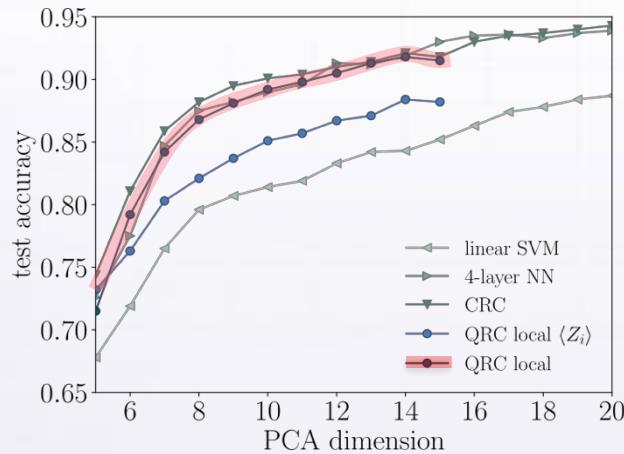
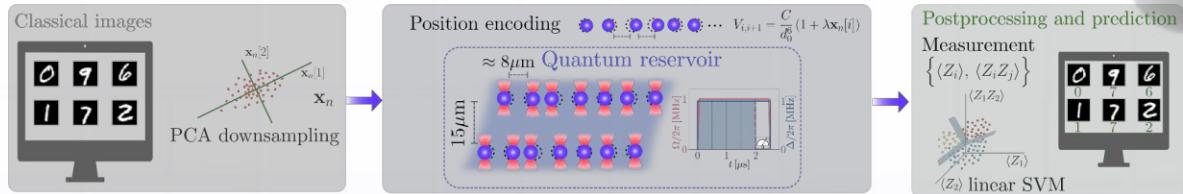
- Trainable
- Resource frugal
- Noise-resilient
- Scalable
- General purpose
- Classically intractable

# Proof of concept – simulations

BLOQADE

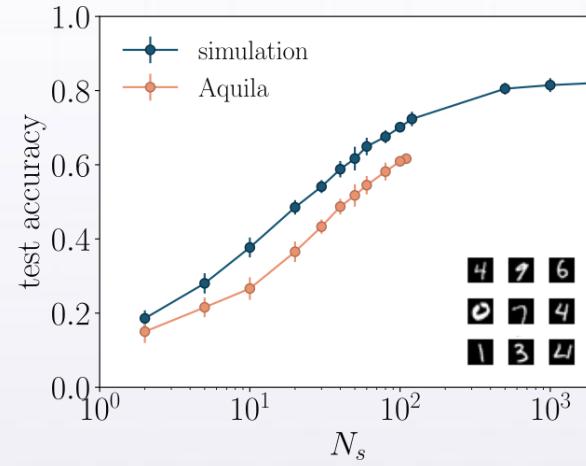
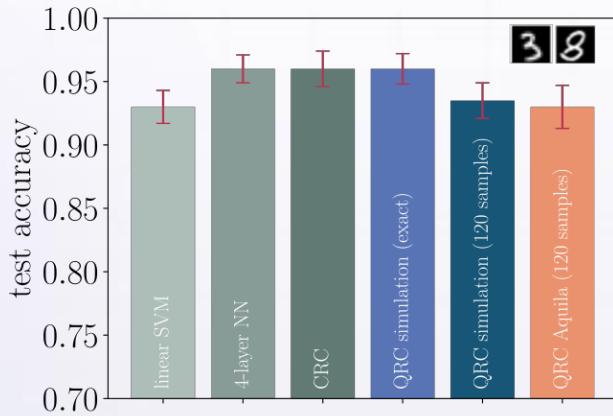
## Performance in simulation on MNIST:

- Exact simulation – reaches problem thresholds
- Practice – sample dependent performance



# Proof of concept – Aquila

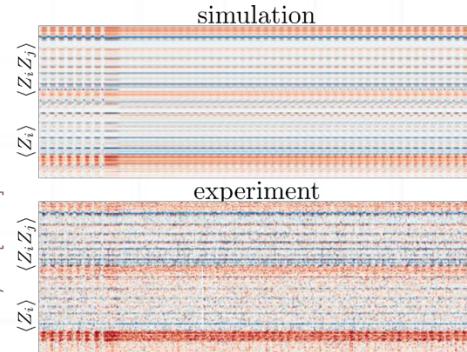
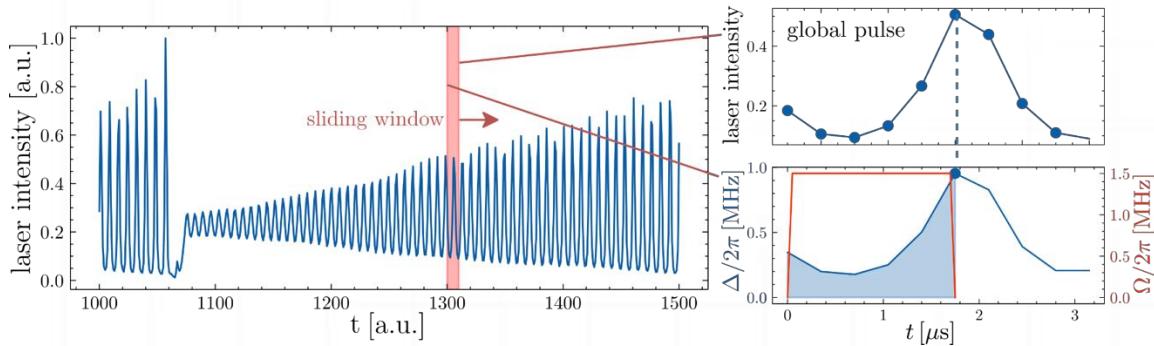
- **Performance on quantum hardware:**
  - 93% test accuracy on Aquila – binary task.
  - Within 0.5% of sampled simulation.
- 62% accuracy – 10-class MNIST.
- Sample requirements higher for multiclass.



# QRC – timeseries and encoding comparison

Look into QRC embeddings:

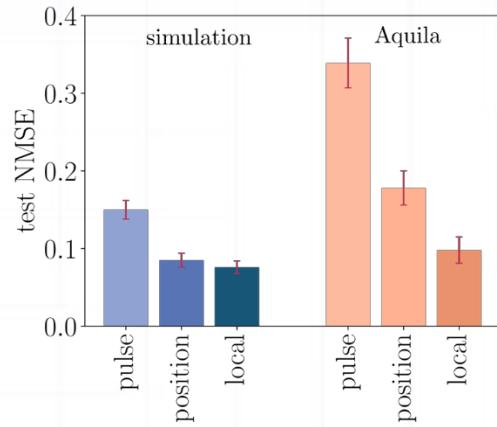
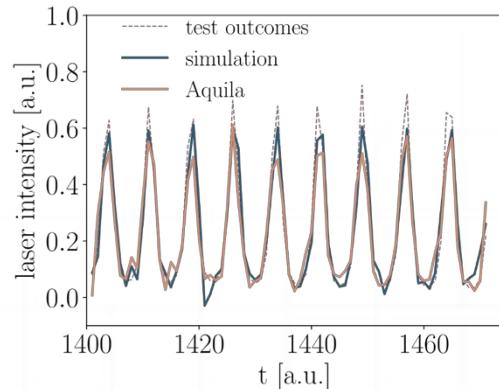
- Embeddings – new expressions of the data.
- Simulation and experiment agree.



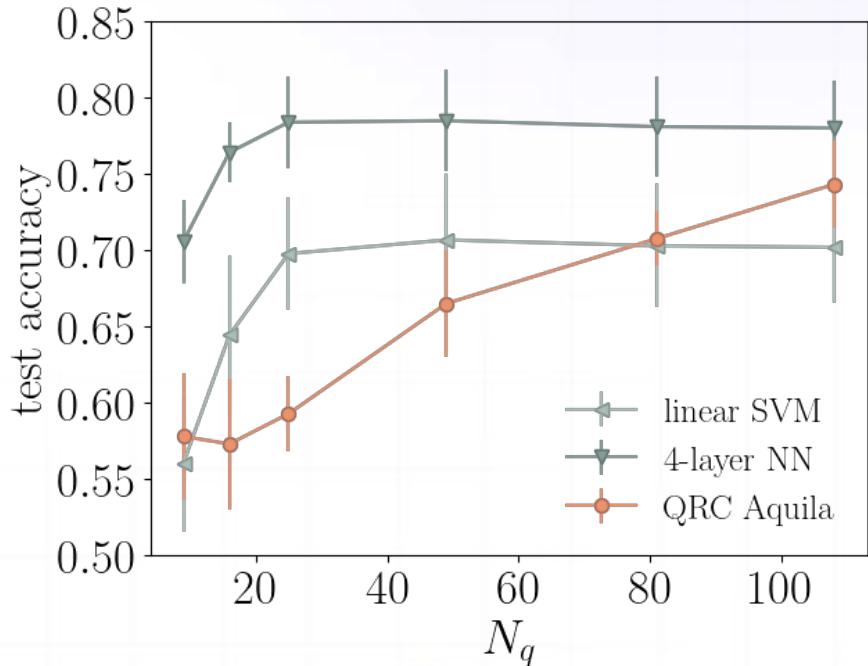
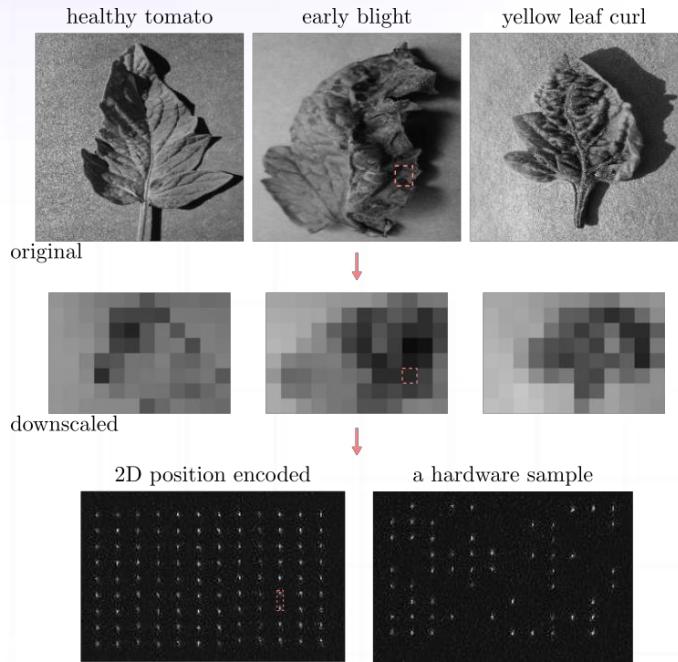
# QRC – timeseries and encoding comparison

## Time series usecase:

- Experiment performing as well as simulation.
- Encodings compared – **local detuning** the best.
- One of the first results with **expanded Aquila capabilities**.



# Scaling up to > 100 qubits – Aquila experiments



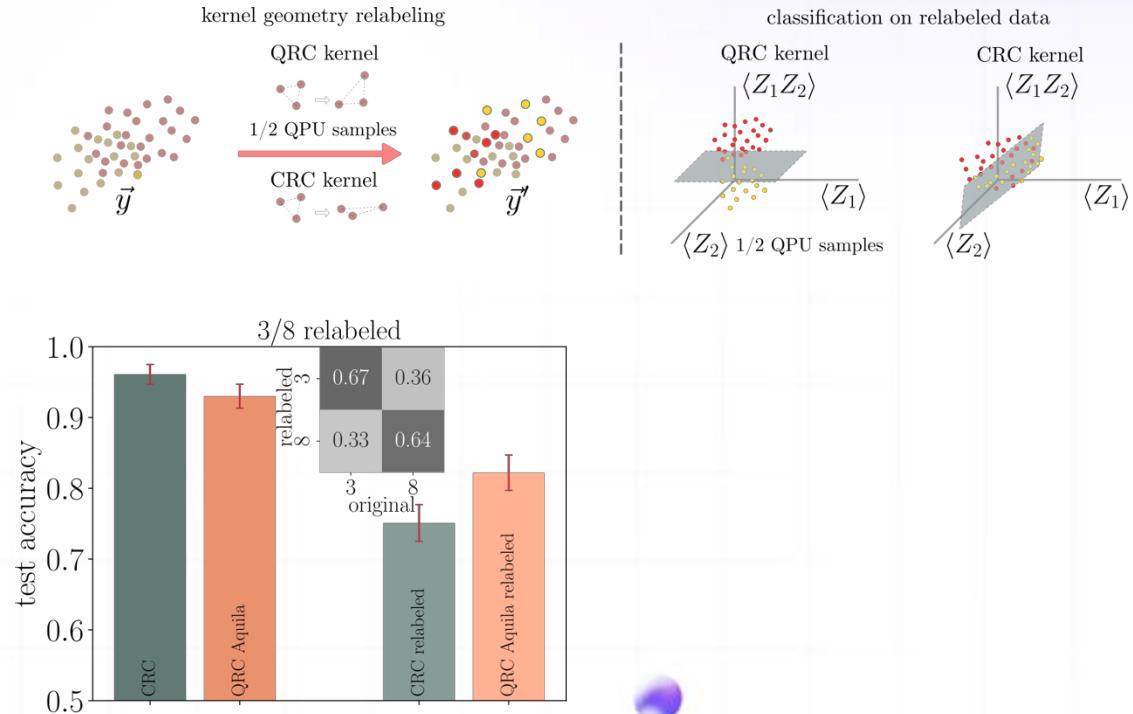
## Main results:

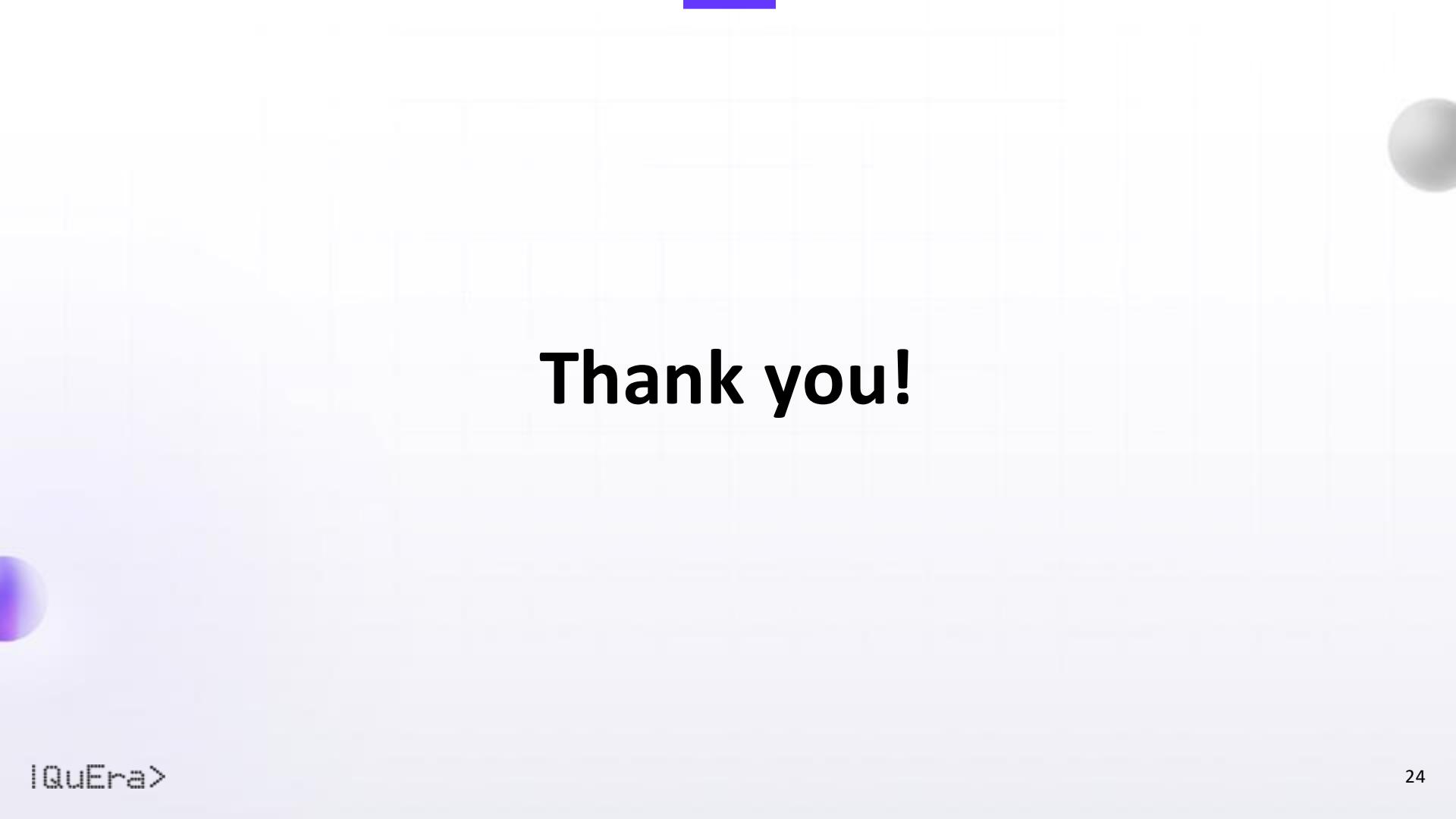
- Consistent scaling with increasing qubit number
- The biggest successful QML experiment to date – 108 qubits

# QRC – Aquila experiments

## Why QRC over classical methods?

- Potential classical intractability
- Data geometry changed in a fundamentally different way
- “Kernel relabeling” construction [Nat. Commun. 12: 2631 (2021)]
- Comparative quantum kernel advantage
- Works on Aquila





# Thank you!