

BENCHMARKING GPU CLUSTERS WITH THE JÜLICH UNIVERSAL QUANTUM COMPUTER SIMULATOR

GTC21 | APRIL 14 | DR. DENNIS WILLSCH



CONTENTS

- 1. Quantum computing
- 2. JUQCS: Simulating quantum computers
- 3. JUQCS-G: Simulating quantum computers on GPUs
- **4. JUQMES:** Simulating physical realizations of quantum computers







Dr. Dennis Willsch



Prof. Dr. Hans De Raedt

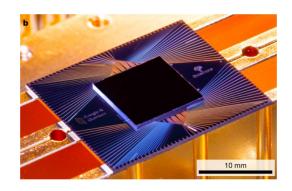


Prof. Dr. Kristel Michielsen



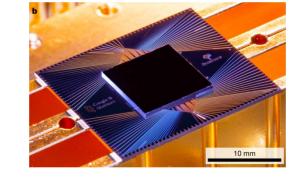
Ideal gate-based quantum computing

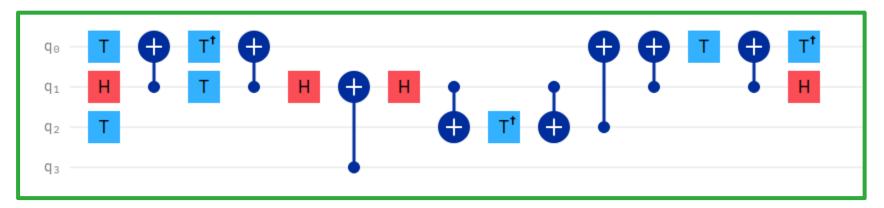
> What does a (gate-based) quantum computer do?





- > What does a (gate-based) quantum computer do?
 - > It runs a quantum circuit

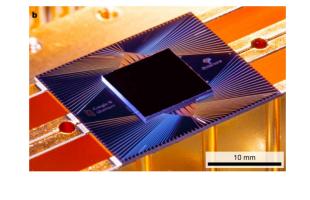


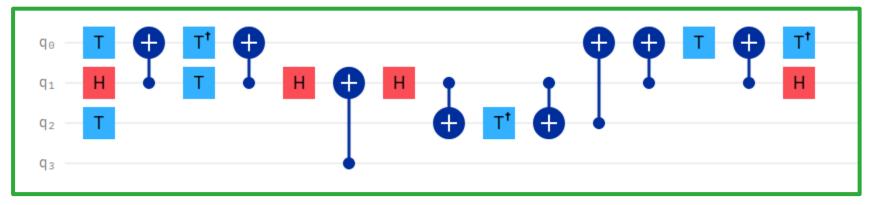




Ideal gate-based quantum computing

- > What does a (gate-based) quantum computer do?
 - > It runs a quantum circuit

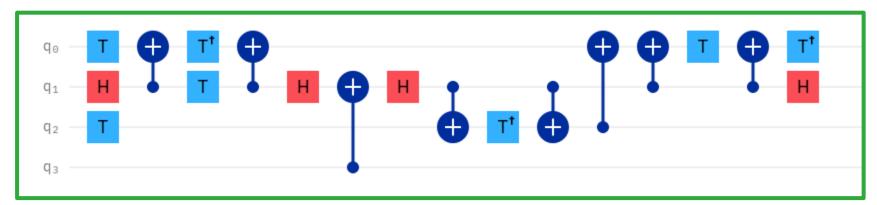




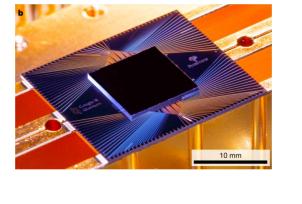
➤ What does this mean, actually?



- > What does a (gate-based) quantum computer do?
 - > It runs a quantum circuit

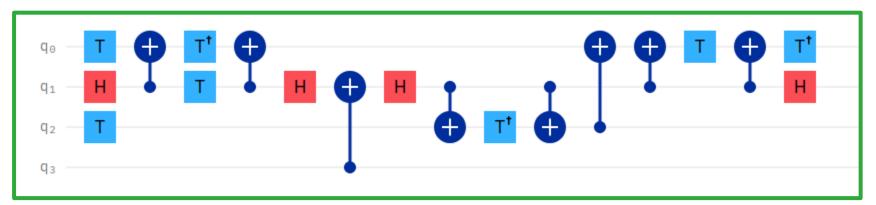


- ➤ What does this mean, actually?
 - > It performs matrix-vector multiplications



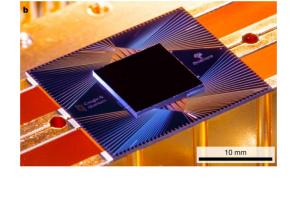


- > What does a (gate-based) quantum computer do?
 - > It runs a quantum circuit



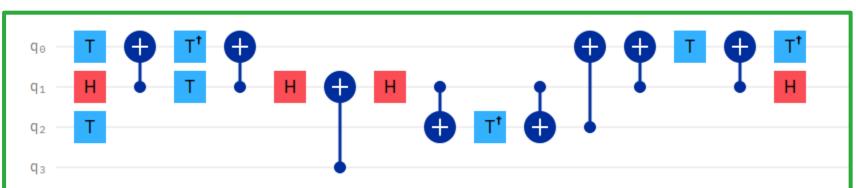
- ➤ What does this mean, actually?
 - > It performs matrix-vector multiplications that are







- > What does a (gate-based) quantum computer do?
 - > It runs a quantum circuit



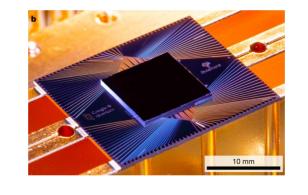
- What does this mean, actually?
 - > It performs matrix-vector multiplications that are

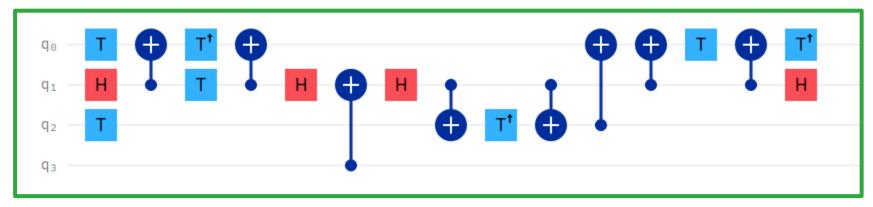






- > What does a (gate-based) quantum computer do?
 - > It runs a quantum circuit





- ➤ What does this mean, actually?
 - > It performs matrix-vector multiplications that are



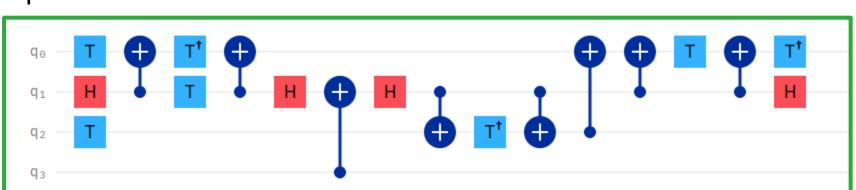






Ideal gate-based quantum computing

- > What does a (gate-based) quantum computer do?
 - > It runs a quantum circuit



- ➤ What does this mean, actually?
 - > It performs matrix-vector multiplications that are



> with huge vectors and huge² matrices



Ideal gate-based quantum computing

> What kind of sparse, unitary matrix-vector multiplications, precisely?



Ideal gate-based quantum computing

vector = state of the QC = 2^n complex numbers

$$|\psi\rangle = \psi_{0\cdots 0}|0\cdots 0\rangle + \cdots + \psi_{1\cdots 1}|1\cdots 1\rangle = \begin{pmatrix} \psi_{0\cdots 0} \\ \vdots \\ \psi_{1\cdots 1} \end{pmatrix}$$

> What kind of sparse, unitary matrix-vector multiplications, precisely?



Ideal gate-based quantum computing

vector = state of the $QC = 2^n$ complex numbers

$$|\psi\rangle = \psi_{0\cdots 0}|0\cdots 0\rangle + \cdots + \psi_{1\cdots 1}|1\cdots 1\rangle = \begin{pmatrix} \psi_{0\cdots 0} \\ \vdots \\ \psi_{1\cdots 1} \end{pmatrix}$$

> What kind of sparse, unitary matrix-vector multiplications, precisely?

each quantum gate = 1 sparse, unitary **matrix**



Ideal gate-based quantum computing

vector = state of the QC = 2^n complex numbers

$$|\psi\rangle = \psi_{0\dots 0}|0\dots 0\rangle + \dots + \psi_{1\dots 1}|1\dots 1\rangle = \begin{pmatrix} \psi_{0\dots 0} \\ \vdots \\ \psi_{1\dots 1} \end{pmatrix}$$

- > What kind of sparse, unitary matrix-vector multiplications, precisely?
 - each quantum gate = 1 sparse, unitary **matrix**

> Example:



Ideal gate-based quantum computing

vector = state of the QC = 2^n complex numbers

$$|\psi\rangle = \psi_{0\cdots 0}|0\cdots 0\rangle + \cdots + \psi_{1\cdots 1}|1\cdots 1\rangle = \begin{pmatrix} \psi_{0\cdots 0} \\ \vdots \\ \psi_{1\cdots 1} \end{pmatrix}$$

> What kind of sparse, unitary matrix-vector multiplications, precisely?

each quantum gate = 1 sparse, unitary **matrix**

> Example:

$$n=4$$
 qubits $2^n=16$ complex numbers



Page 4

Ideal gate-based quantum computing

vector = state of the QC = 2^n complex numbers

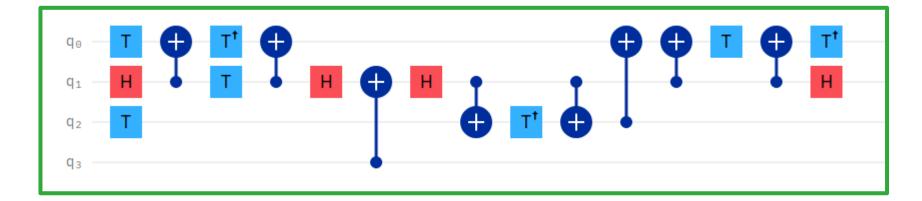
$$|\psi\rangle = \psi_{0\cdots 0}|0\cdots 0\rangle + \cdots + \psi_{1\cdots 1}|1\cdots 1\rangle = \begin{pmatrix} \psi_{0\cdots 0} \\ \vdots \\ \psi_{1\cdots 1} \end{pmatrix}$$

> What kind of sparse, unitary matrix-vector multiplications, precisely?

each quantum gate = 1 sparse, unitary **matrix**

> Example:

$$n = 4$$
 qubits $2^n = 16$ complex numbers



Ideal gate-based quantum computing

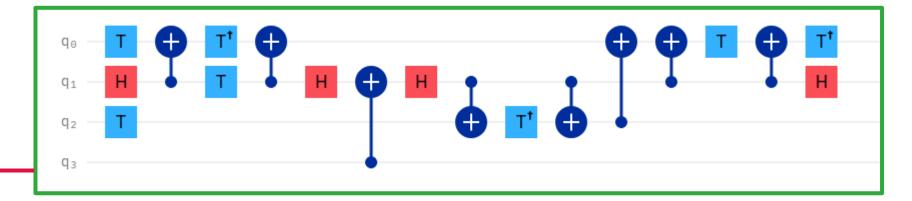
vector = state of the QC = 2^n complex numbers

$$|\psi\rangle = \psi_{0\cdots 0}|0\cdots 0\rangle + \cdots + \psi_{1\cdots 1}|1\cdots 1\rangle = \begin{pmatrix} \psi_{0\cdots 0} \\ \vdots \\ \psi_{1\cdots 1} \end{pmatrix}$$

- > What kind of sparse, unitary matrix-vector multiplications, precisely?
 - each quantum gate = 1 sparse, unitary **matrix**

> Example:

$$n = 4$$
 qubits $2^n = 16$ complex numbers



Initial state of QC:

$$|\psi\rangle = |0000\rangle = \begin{pmatrix} 1\\0\\\vdots\\0 \end{pmatrix}$$

Page 4

Ideal gate-based quantum computing

vector = state of the QC = 2^n complex numbers

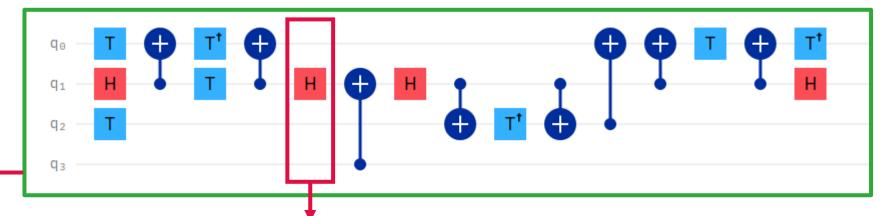
$$|\psi\rangle = \psi_{0\cdots 0}|0\cdots 0\rangle + \cdots + \psi_{1\cdots 1}|1\cdots 1\rangle = \begin{pmatrix} \psi_{0\cdots 0} \\ \vdots \\ \psi_{1\cdots 1} \end{pmatrix}$$

> What kind of sparse, unitary matrix-vector multiplications, precisely?

each quantum gate = 1 sparse, unitary **matrix**

> Example:

$$n = 4$$
 qubits $2^n = 16$ complex numbers



Initial state of QC:

$$|\psi\rangle = |0000\rangle = \begin{pmatrix} 1\\0\\\vdots\\0 \end{pmatrix}$$

Example matrix-vector multiplication: \blacksquare on qubit q_1

Ideal gate-based quantum computing

vector = state of the QC = 2^n complex numbers

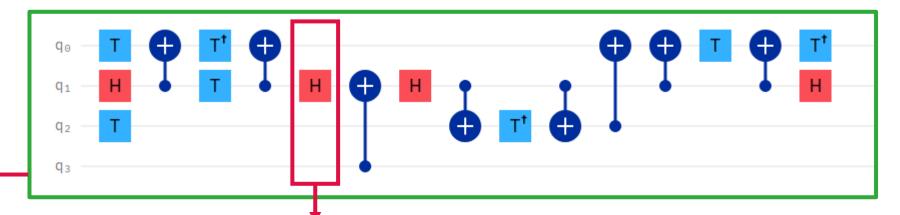
$$|\psi\rangle = \psi_{0\cdots 0}|0\cdots 0\rangle + \cdots + \psi_{1\cdots 1}|1\cdots 1\rangle = \begin{pmatrix} \psi_{0\cdots 0} \\ \vdots \\ \psi_{1\cdots 1} \end{pmatrix}$$

> What kind of sparse, unitary matrix-vector multiplications, precisely?

each quantum gate = 1 sparse, unitary **matrix**

> Example:

$$n = 4$$
 qubits $2^n = 16$ complex numbers



Initial state of QC:

$$|\psi\rangle = |0000\rangle = \begin{pmatrix} 1\\0\\\vdots\\0 \end{pmatrix}$$

Example matrix-vector multiplication: \blacksquare on qubit q_1

Page 4

For each q_3 , q_2 , q_0 perform 2x2 update:

Ideal gate-based quantum computing

vector = state of the QC = 2^n complex numbers

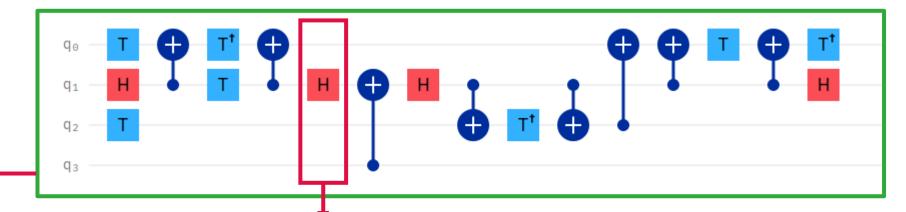
$$|\psi\rangle = \psi_{0\cdots 0}|0\cdots 0\rangle + \cdots + \psi_{1\cdots 1}|1\cdots 1\rangle = \begin{pmatrix} \psi_{0\cdots 0} \\ \vdots \\ \psi_{1\cdots 1} \end{pmatrix}$$

> What kind of sparse, unitary matrix-vector multiplications, precisely?

each quantum gate = 1 sparse, unitary **matrix**

> Example:

$$n = 4$$
 qubits $2^n = 16$ complex numbers



Initial state of QC:

$$|\psi\rangle = |0000\rangle = \begin{pmatrix} 1\\0\\ \vdots\\0 \end{pmatrix}$$

Example matrix-vector multiplication: \blacksquare on qubit q_1

For each
$$q_3$$
, q_2 , q_0 perform 2x2 update: $\begin{pmatrix} \psi_{q_3q_20q_0} \\ \psi_{q_3q_21q_0} \end{pmatrix} \leftarrow \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} \psi_{q_3q_20q_0} \\ \psi_{q_3q_21q_0} \end{pmatrix}$

Page 4

Ideal gate-based quantum computing

vector = state of the QC = 2^n complex numbers

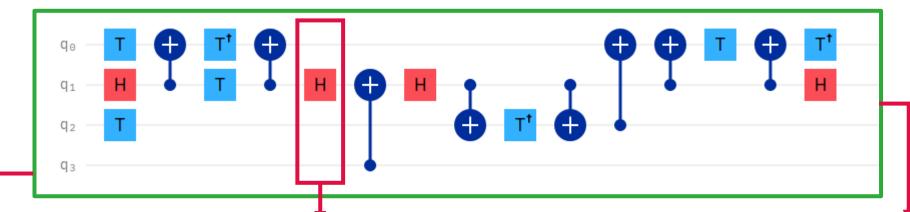
$$|\psi\rangle = \psi_{0\cdots 0}|0\cdots 0\rangle + \cdots + \psi_{1\cdots 1}|1\cdots 1\rangle = \begin{pmatrix} \psi_{0\cdots 0} \\ \vdots \\ \psi_{1\cdots 1} \end{pmatrix}$$

> What kind of sparse, unitary matrix-vector multiplications, precisely?

each quantum gate = 1 sparse, unitary **matrix**

> Example:

$$n = 4$$
 qubits $2^n = 16$ complex numbers



Initial state of QC:

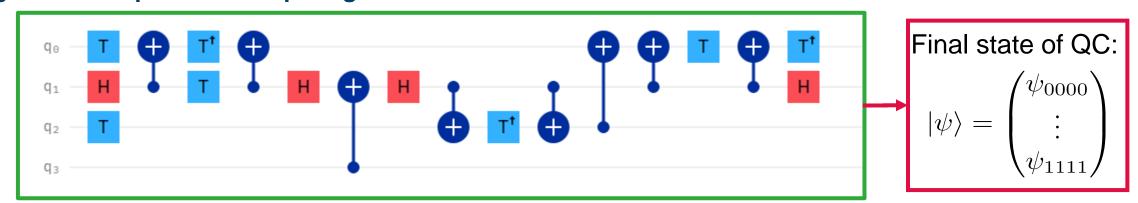
$$|\psi\rangle = |0000\rangle = \begin{pmatrix} 1\\0\\\vdots\\0 \end{pmatrix}$$

Example matrix-vector multiplication: \blacksquare on qubit q_1

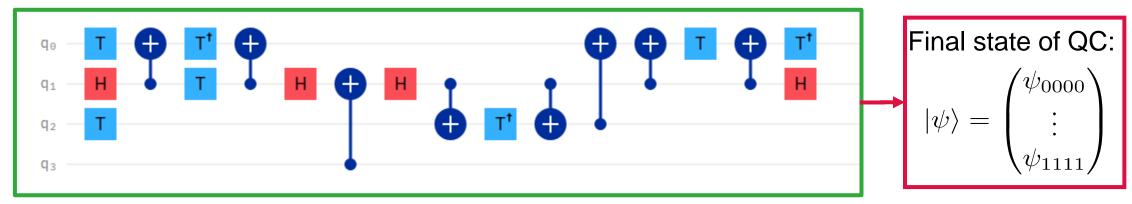
For each
$$q_3$$
, q_2 , q_0 perform 2x2 update: $\begin{pmatrix} \psi_{q_3q_20q_0} \\ \psi_{q_3q_21q_0} \end{pmatrix} \leftarrow \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & 1 \\ 1 & -1 \end{pmatrix} \begin{pmatrix} \psi_{q_3q_20q_0} \\ \psi_{q_3q_21q_0} \end{pmatrix} \hspace{0.2cm} |\psi\rangle = \begin{pmatrix} |\psi\rangle| + |\psi\rangle|\psi$

Final state of QC:

$$|\psi\rangle = \begin{pmatrix} \psi_{0000} \\ \vdots \\ \psi_{1111} \end{pmatrix}$$

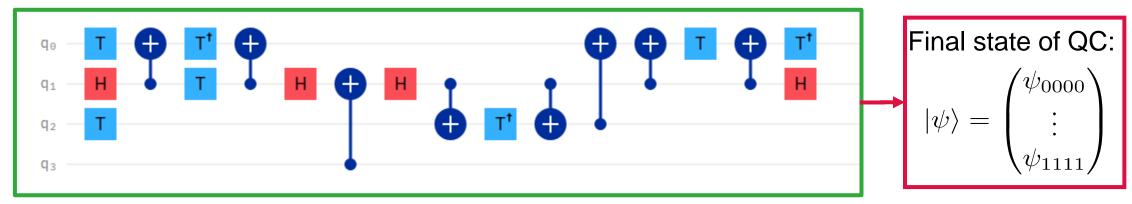


Ideal gate-based quantum computing



> What does a hardware realization of a QC return?

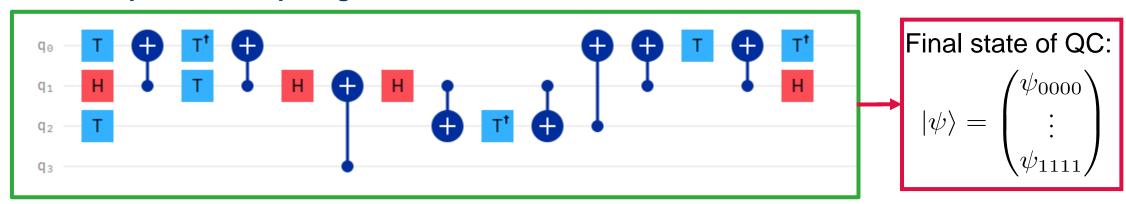




- > What does a **hardware realization** of a QC return?
 - > The quantum state after all sparse matrix-vector multiplications?



Ideal gate-based quantum computing

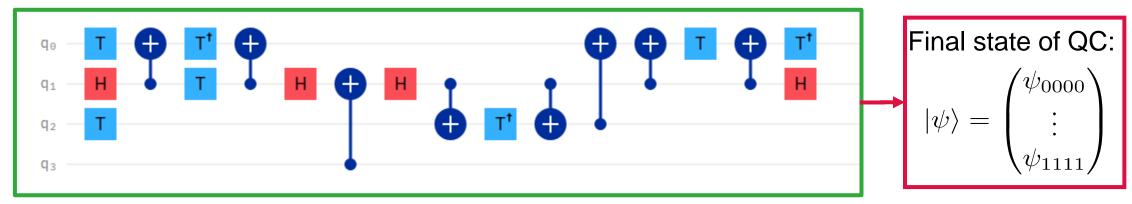


- > What does a hardware realization of a QC return?
 - > The quantum state after all sparse matrix-vector multiplications?
 - \triangleright No! That would be 2^n complex numbers.

For 40 qubits: $2^{40} \psi' s$ = 16 TiB complex numbers



Ideal gate-based quantum computing



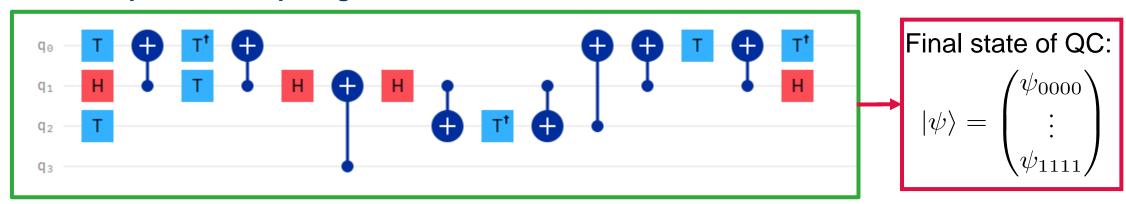
- > What does a hardware realization of a QC return?
 - > The quantum state after all sparse matrix-vector multiplications?
 - \triangleright No! That would be 2^n complex numbers.

For 40 qubits: $2^{40} \, \psi' \mathrm{s}$ = 16 TiB complex numbers

 \triangleright What then? Only a single bitstring $q_{n-1} \cdots q_1 q_0$ with n bits



Ideal gate-based quantum computing



- What does a hardware realization of a QC return?
 - > The quantum state after all sparse matrix-vector multiplications?
 - \triangleright No! That would be 2^n complex numbers.

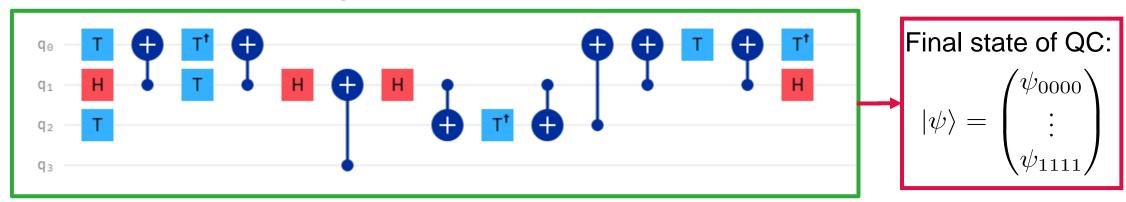
For 40 qubits: $2^{40} \psi' s = 16$ TiB complex numbers

- \triangleright What then? Only a single bitstring $q_{n-1} \cdots q_1 q_0$ with n bits
- > The complex numbers only define the **probability**:

$$|\psi_{q_{n-1}\cdots q_1q_0}|^2$$
 = probability to return bitstring $q_{n-1}\cdots q_1q_0$

Page 5

Ideal gate-based quantum computing



- > What does a hardware realization of a QC return?
 - > The quantum state after all sparse matrix-vector multiplications?
 - \triangleright No! That would be 2^n complex numbers.

For 40 qubits: $2^{40} \, \psi' \mathrm{s}$ = 16 TiB complex numbers

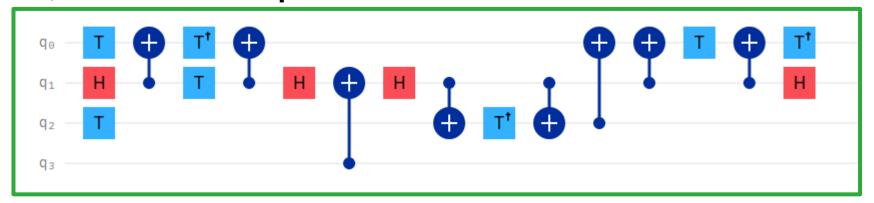
- \triangleright What then? Only a single bitstring $q_{n-1} \cdots q_1 q_0$ with n bits
- > The complex numbers only define the **probability**:

$$|\psi_{q_{n-1}\cdots q_1q_0}|^2$$
 = probability to return bitstring $q_{n-1}\cdots q_1q_0$

> Need to run circuit repeatedly to sample from the distribution

Ideal gate-based quantum computing

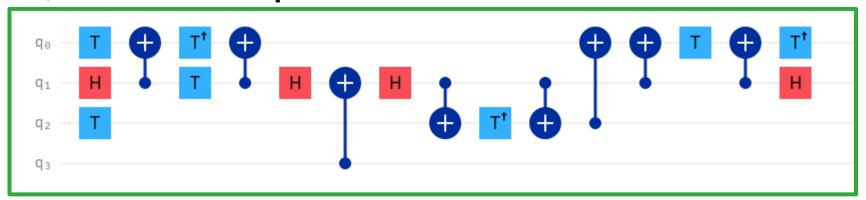
➤ In particular, what does this quantum circuit do?





Ideal gate-based quantum computing

➤ In particular, what does this quantum circuit do?



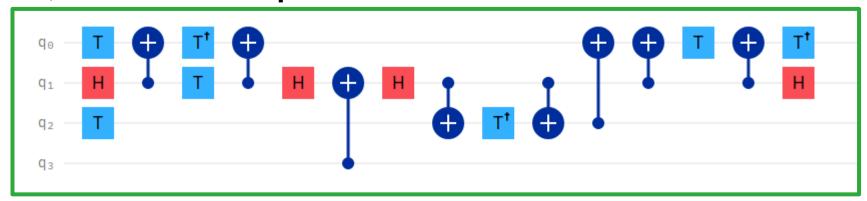
➤ 2-qubit adder

$$|q_3q_2\rangle|q_1q_0\rangle \mapsto |q_3q_2\rangle|q_3q_2 + q_1q_0\rangle$$



Ideal gate-based quantum computing

> In particular, what does this quantum circuit do?

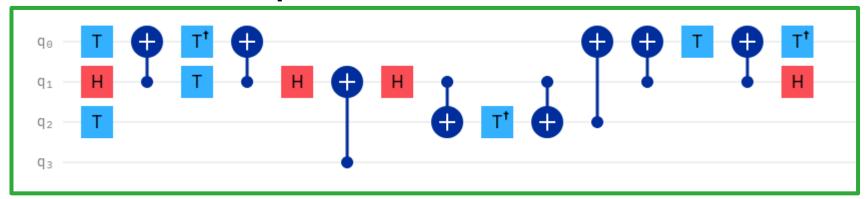


- ➤ 2-qubit adder
 - ightharpoonup e.g. $|2\rangle|1\rangle\mapsto|2\rangle|3\rangle$

$$|q_3q_2\rangle|q_1q_0\rangle \mapsto |q_3q_2\rangle|q_3q_2 + q_1q_0\rangle$$

Ideal gate-based quantum computing

> In particular, what does this quantum circuit do?



➤ 2-qubit adder

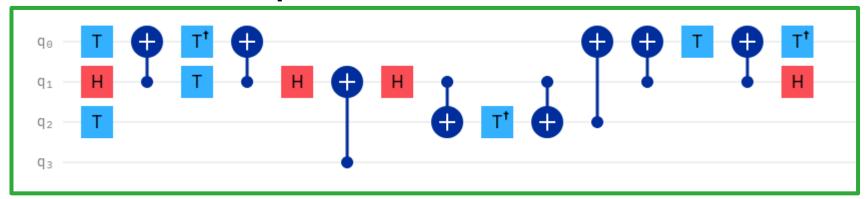
$$|q_3q_2\rangle|q_1q_0\rangle \mapsto |q_3q_2\rangle|q_3q_2 + q_1q_0\rangle$$

Page 6

- ightharpoonup e.g. $|2\rangle|1\rangle\mapsto|2\rangle|3\rangle$
- > but also **superpositions**:

Ideal gate-based quantum computing

> In particular, what does this quantum circuit do?



➤ 2-qubit adder

$$|q_3q_2\rangle|q_1q_0\rangle \mapsto |q_3q_2\rangle|q_3q_2 + q_1q_0\rangle$$

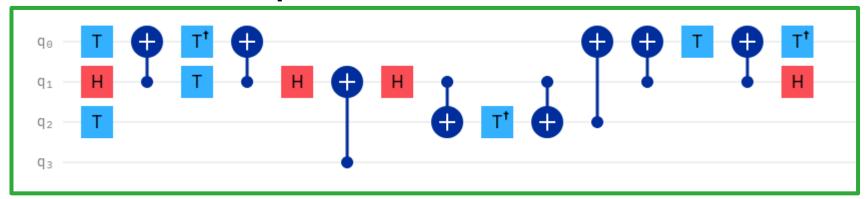
- ightharpoonup e.g. $|2\rangle|1\rangle\mapsto|2\rangle|3\rangle$
- but also superpositions:

$$|2\rangle \frac{|0\rangle + |1\rangle + |2\rangle}{\sqrt{3}} \mapsto |2\rangle \frac{|2\rangle + |3\rangle + |0\rangle}{\sqrt{3}}$$



Ideal gate-based quantum computing

> In particular, what does this quantum circuit do?



➤ 2-qubit adder

$$|q_3q_2\rangle|q_1q_0\rangle \mapsto |q_3q_2\rangle|q_3q_2 + q_1q_0\rangle$$

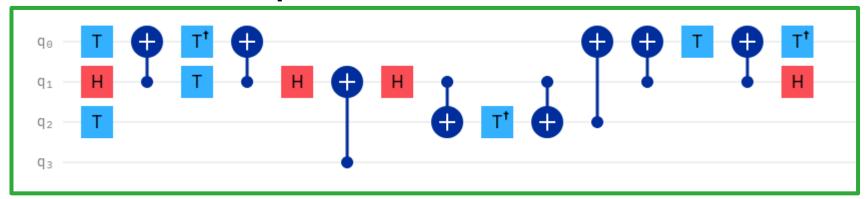
- ightharpoonup e.g. $|2\rangle|1\rangle\mapsto|2\rangle|3\rangle$
- > but also **superpositions**:

$$|2\rangle \frac{|0\rangle + |1\rangle + |2\rangle}{\sqrt{3}} \mapsto |2\rangle \frac{|2\rangle + |3\rangle + |0\rangle}{\sqrt{3}}$$

$$\begin{pmatrix} \vdots \\ \psi_{1000} = 1/\sqrt{3} \\ \psi_{1001} = 1/\sqrt{3} \\ \psi_{1010} = 1/\sqrt{3} \\ \psi_{1011} = 0 \\ \vdots \end{pmatrix} \mapsto \begin{pmatrix} \vdots \\ \psi_{1000} = 1/\sqrt{3} \\ \psi_{1001} = 0 \\ \psi_{1010} = 1/\sqrt{3} \\ \psi_{1011} = 1/\sqrt{3} \\ \vdots \end{pmatrix}$$

Ideal gate-based quantum computing

> In particular, what does this quantum circuit do?



> 2-qubit adder

$$|q_3q_2\rangle|q_1q_0\rangle \mapsto |q_3q_2\rangle|q_3q_2 + q_1q_0\rangle$$

- ightharpoonup e.g. $|2\rangle|1\rangle\mapsto|2\rangle|3\rangle$
- > but also **superpositions**:

$$|2\rangle \frac{|0\rangle + |1\rangle + |2\rangle}{\sqrt{3}} \mapsto |2\rangle \frac{|2\rangle + |3\rangle + |0\rangle}{\sqrt{3}}$$

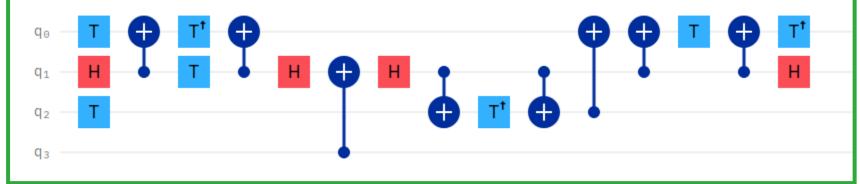
$$\begin{pmatrix} \vdots \\ \psi_{1000} = 1/\sqrt{3} \\ \psi_{1001} = 1/\sqrt{3} \\ \psi_{1010} = 1/\sqrt{3} \\ \psi_{1011} = 0 \\ \vdots \end{pmatrix} \mapsto \begin{pmatrix} \vdots \\ \psi_{1000} = 1/\sqrt{3} \\ \psi_{1001} = 0 \\ \psi_{1010} = 1/\sqrt{3} \\ \psi_{1011} = 1/\sqrt{3} \\ \vdots \end{pmatrix}$$

JUQCS

Jülich universal quantum computer simulator

- > What does a quantum computer **simulator** do?
 - > It runs a quantum circuit



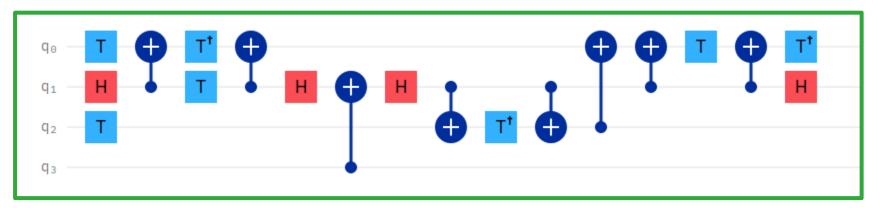




Jülich universal quantum computer simulator

- What does a quantum computer simulator do?
 - > It runs a quantum circuit





Page 7

- ➤ What does this mean, actually?
 - > It performs matrix-vector multiplications that are



> with huge vectors and huge² matrices



Distribution of the quantum state

How does the simulator manage all these complex numbers?



Distribution of the quantum state

How does the simulator manage all these complex numbers?

 \rightarrow Distribute quantum state $|\psi\rangle=(\psi_{\cdots q_2q_1q_0})$ over multiple compute nodes

Distribution of the quantum state

How does the simulator manage all these complex numbers?

 \rightarrow Distribute quantum state $|\psi\rangle=(\psi_{\cdots q_2q_1q_0})$ over multiple compute nodes

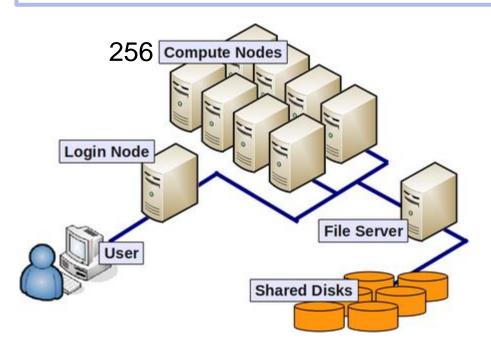
For 40 qubits: $2^{40} \psi's = 16$ TiB complex numbers = 64 GiB per node with 256 nodes

Distribution of the quantum state

How does the simulator manage all these complex numbers?

 \rightarrow Distribute quantum state $|\psi\rangle=(\psi_{\cdots q_2q_1q_0})$ over multiple compute nodes

For 40 qubits: $2^{40} \psi'_s = 16$ TiB complex numbers = 64 GiB per node with 256 nodes



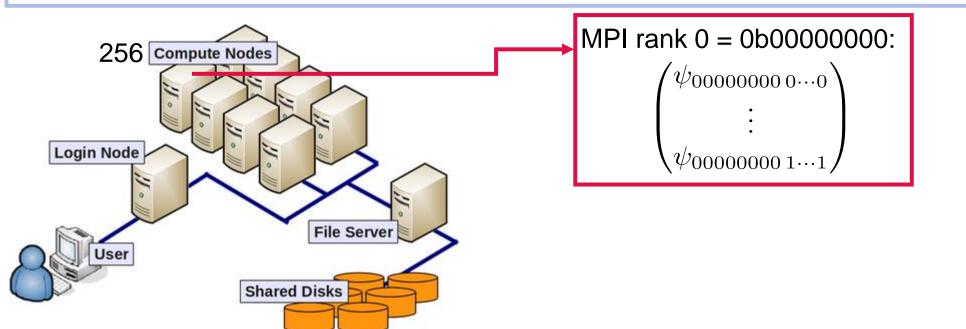


Distribution of the quantum state

How does the simulator manage all these complex numbers?

 \rightarrow Distribute quantum state $|\psi\rangle=(\psi_{\cdots q_2q_1q_0})$ over multiple compute nodes

For 40 qubits: $2^{40} \psi' s = 16$ TiB complex numbers = 64 GiB per node with 256 nodes

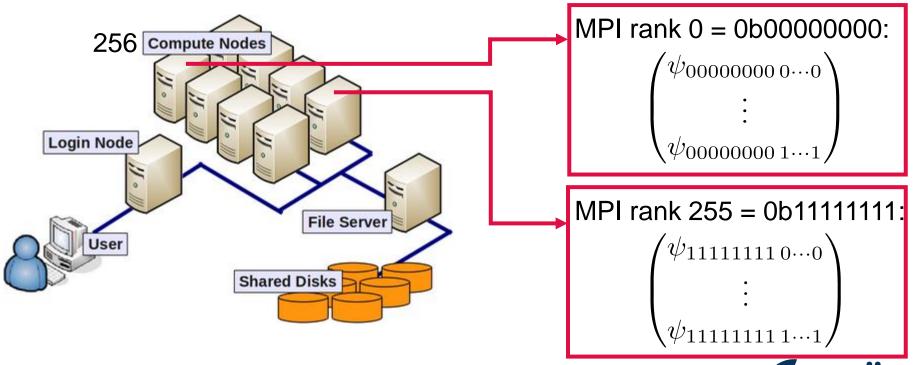


Distribution of the quantum state

How does the simulator manage all these complex numbers?

 \rightarrow Distribute quantum state $|\psi\rangle=(\psi_{\cdots q_2q_1q_0})$ over multiple compute nodes

For 40 qubits: $2^{40} \psi's = 16$ TiB complex numbers = 64 GiB per node with 256 nodes

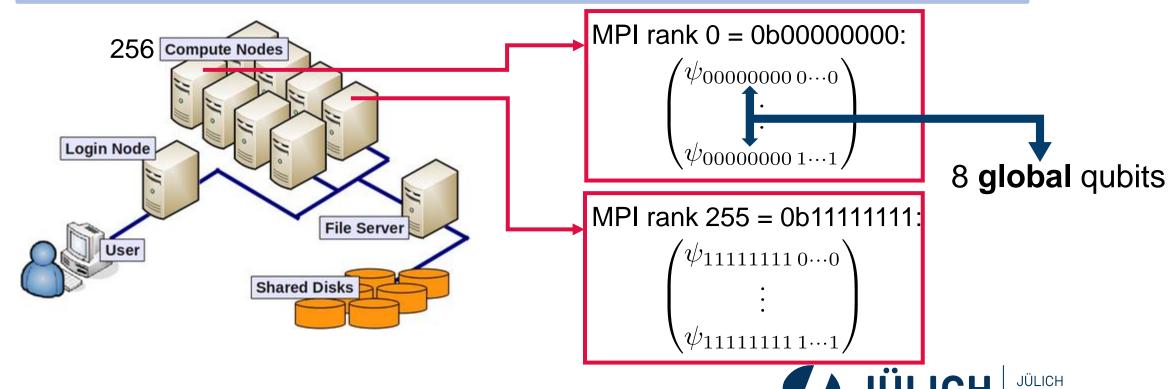


Distribution of the quantum state

How does the simulator manage all these complex numbers?

 \rightarrow Distribute quantum state $|\psi\rangle=(\psi_{\cdots q_2q_1q_0})$ over multiple compute nodes

For 40 qubits: $2^{40} \psi's = 16$ TiB complex numbers = 64 GiB per node with 256 nodes



SUPERCOMPUTING

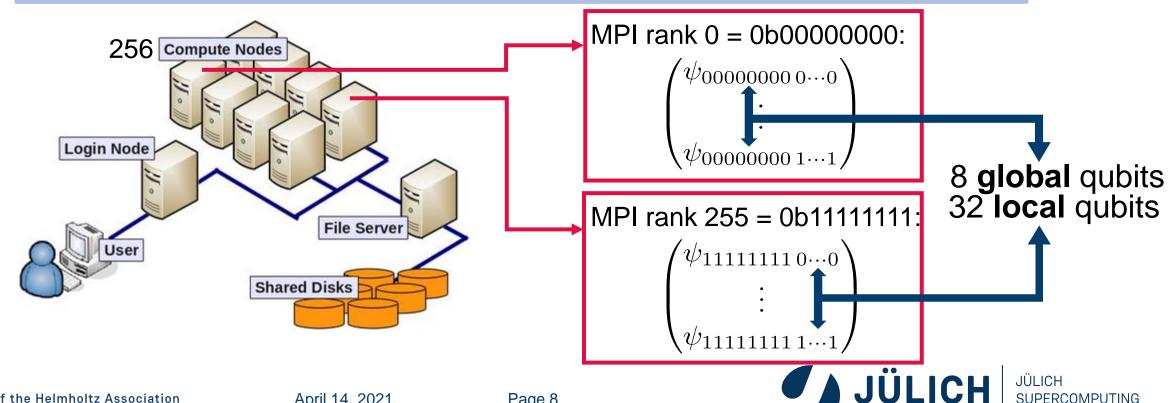
CENTRE

Distribution of the quantum state

How does the simulator manage all these complex numbers?

 \rightarrow Distribute quantum state $|\psi\rangle = (\psi ... q_2 q_1 q_0)$ over multiple compute nodes

For 40 qubits: $2^{40} \psi'_s = 16$ TiB complex numbers = 64 GiB per node with 256 nodes



CENTRE

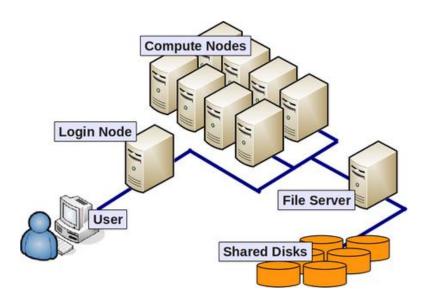
Simulating quantum computers on GPUs CUDA MPI Fortran

➤ Distribute quantum state on GPUs (NVIDIA A100: 40GB per GPU)



Simulating quantum computers on GPUs CUDA MPI Fortran

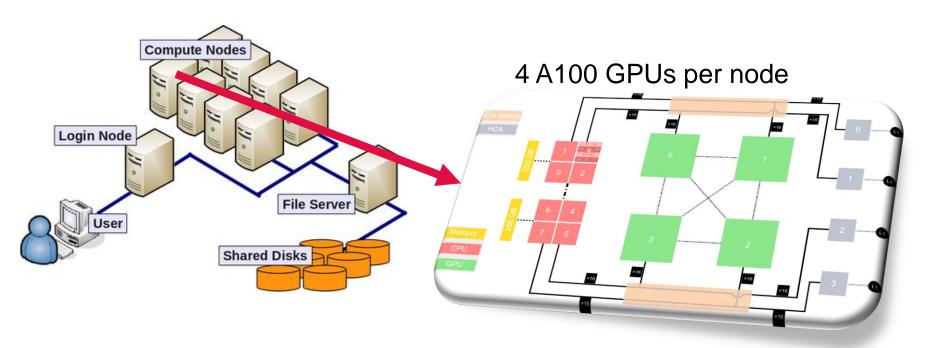
➤ Distribute quantum state on GPUs (NVIDIA A100: 40GB per GPU)





Simulating quantum computers on GPUs CUDA MPI Fortran

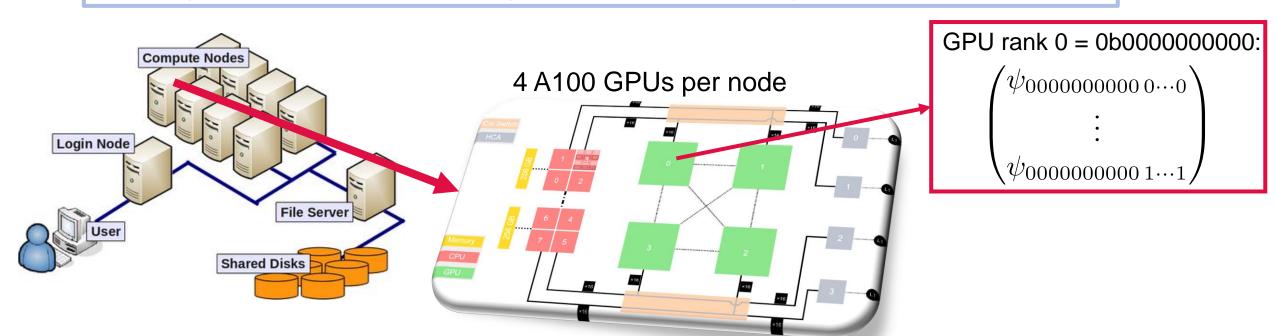
➤ Distribute quantum state on GPUs (NVIDIA A100: 40GB per GPU)





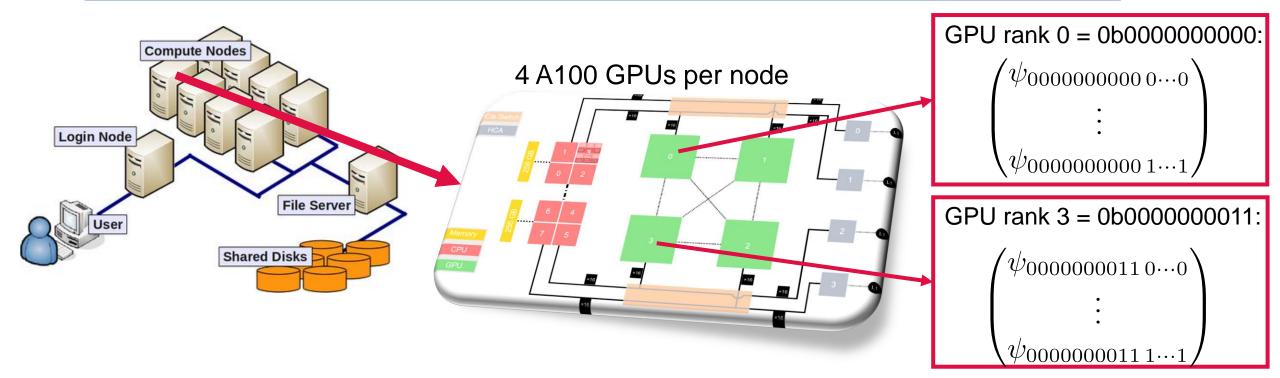
Simulating quantum computers on GPUs CUDA MPI Fortran

➤ Distribute quantum state on GPUs (NVIDIA A100: 40GB per GPU)



Simulating quantum computers on GPUs CUDA MPI Fortran

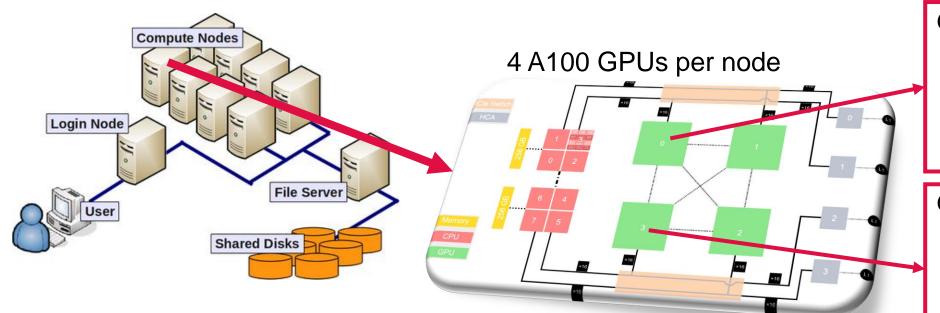
➤ Distribute quantum state on GPUs (NVIDIA A100: 40GB per GPU)



Simulating quantum computers on GPUs CUDA MPI Fortran

> Distribute quantum state on GPUs (NVIDIA A100: 40GB per GPU)

For 40 qubits: $2^{40} \psi' s = 16$ TiB complex numbers = 16 GiB per GPU with 4*256 GPUs



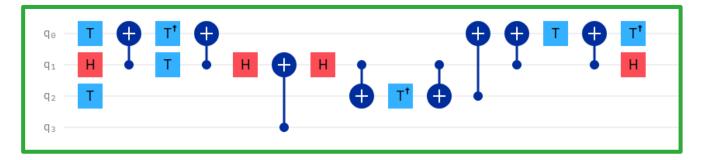
GPU rank 0 = 0b00000000000: $\psi_{000000000001\cdots 1}$.

GPU rank 3 = 0b0000000011: $(\psi_{0000000011\,0...0})$ $\psi_{0000000011\,1\cdots 1}$

> The MPI communication scheme is the same

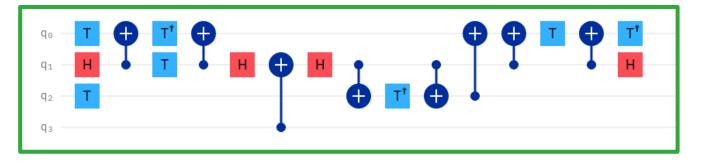


CUDA implementation





CUDA implementation

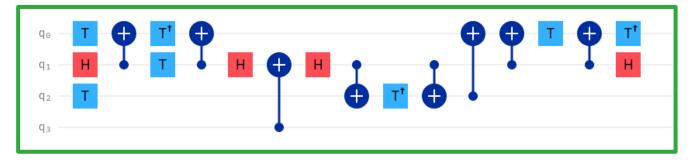


How to implement these matrix-vector multiplications in the most efficient way?

```
\psi_{000000000000000000}
\psi_{0000000001...1}
\psi_{000000010...0}
\psi_{000000011...1}
\psi_{1111111110\cdots 0}
```



CUDA implementation

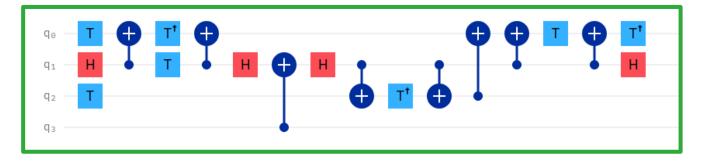


How to implement these matrix-vector multiplications in the most efficient way?

```
\psi000000000 0\cdots0
  GPU rank 0
\psi_{0000000001\cdots 1}
\psi_{000000010...0}
\psi_{000000011...1}
\psi_{1111111110\cdots 0}
```



CUDA implementation

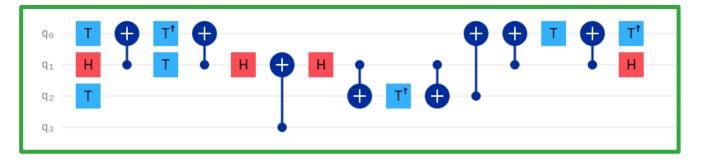


How to implement these matrix-vector multiplications in the most efficient way?

```
\psi_{000000000000\cdots0}
   GPU rank 0
\psi_{0000000001...1}
\psi_{00000001\,0\cdots0}
  GPU rank 1
\psi_{00000001\,1\cdots 1}
\psi_{1111111110\cdots 0}
```



CUDA implementation

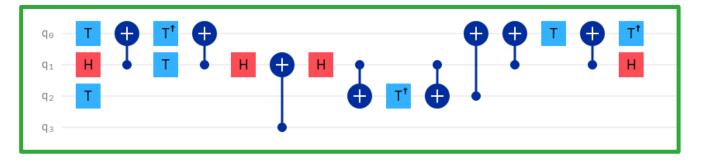


How to implement these matrix-vector multiplications in the most efficient way?

```
\psi_{000000000000\cdots0}
  GPU rank 0
\psi_{0000000001\cdots 1}
\psi_{00000001\,0\cdots0}
  GPU rank 1
\psi_{00000001\,1\cdots 1}
\psi_{1111111110\cdots 0}
GPU rank 255
```



CUDA implementation



How to implement these matrix-vector multiplications in the most efficient way?

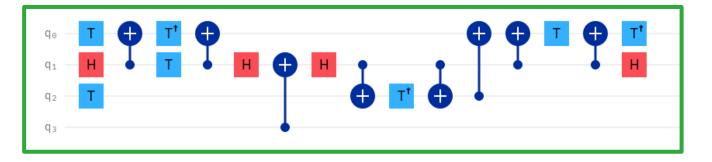
Full quantum state:

```
\psi_{000000000000\cdots0}
  GPU rank 0
\psi_{0000000001\cdots 1}
\psi_{00000001\,0\cdots0}
  GPU rank 1
\psi_{00000001\,1\cdots 1}
\psi_{1111111110\cdots 0}
GPU rank 255
```

Quantum gate on **local** qubits:

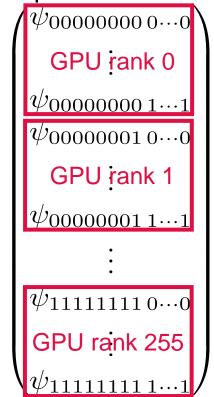


CUDA implementation



How to implement these matrix-vector multiplications in the most efficient way?

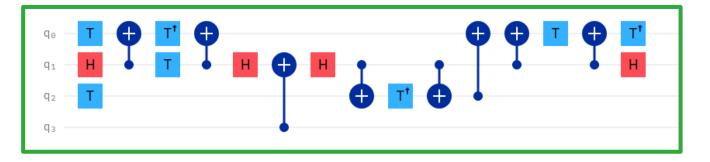
Full quantum state:



Quantum gate on **local** qubits: on qubit q_{30}

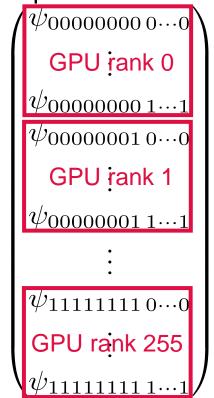


CUDA implementation



How to implement these matrix-vector multiplications in the most efficient way?

Full quantum state:



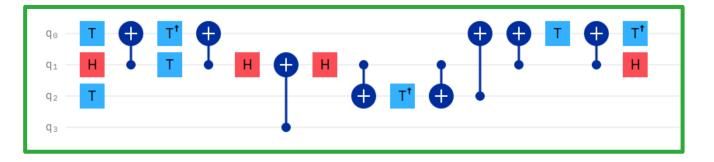
Quantum gate on **local** qubits:

e.g. H on qubit q_{30}

→ Each MPI rank r performs local 2x2 updates of the form

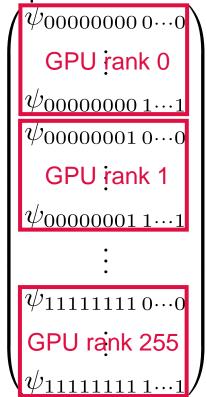


CUDA implementation



How to implement these matrix-vector multiplications in the most efficient way?

Full quantum state:



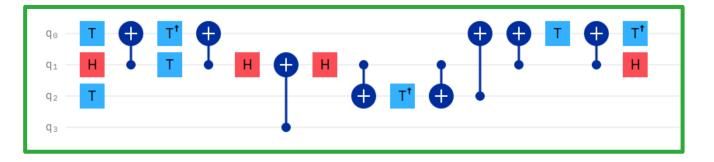
Quantum gate on **local** qubits:

e.g. H on qubit
$$q_{30}$$

 \rightarrow Each MPI rank r performs local 2x2 updates of the form

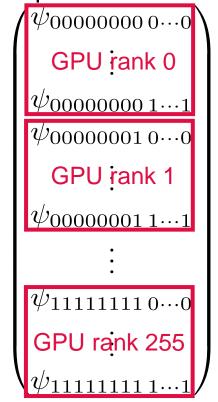
$$\begin{pmatrix} \psi_{rrrrrrr*0*\cdots*} \\ \psi_{rrrrrr*1*\cdots*} \end{pmatrix} \leftarrow \mathbf{H} \quad \begin{pmatrix} \psi_{rrrrrrr*0*\cdots*} \\ \psi_{rrrrrrr*1*\cdots*} \end{pmatrix}$$

CUDA implementation



How to implement these matrix-vector multiplications in the most efficient way?

Full quantum state:

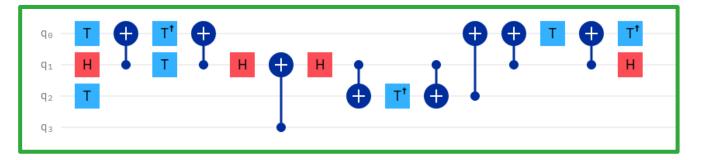


Quantum gate on **local** qubits:

e.g. H on qubit
$$q_{30}$$

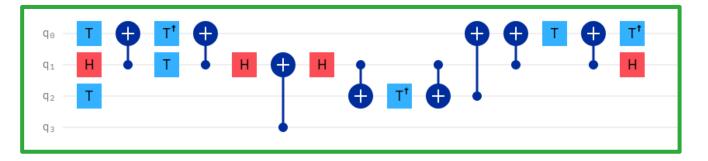
 \rightarrow Each MPI rank r performs local 2x2 updates of the form

CUDA implementation



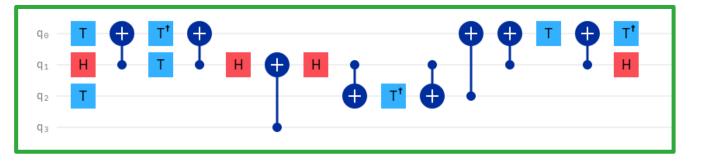


CUDA implementation



```
attributes(global) subroutine H operation GPU(nstates,psi R,psi I,i0,i1,c)
        USE cudafor
        implicit real(kind=8) (a-h,o-z)
        integer(kind=8),parameter :: one=1
        integer(kind=8), value :: nstates,i1
        integer(kind=8) :: m1,nm1,i,j,k
        integer(kind=4), value :: i0
        real(kind=8), value :: c
        real(kind=8),dimension(0:nstates-1),device:: psi R,psi I
        k=blockidx%x-1
        k=k*blockdim%x+threadidx%x-1
        m1=i1-1
        nm1=not(m1)
        i=ishft(iand(k,nm1),one)+iand(k,m1)
        j=i+i1
        r0=psi R(i)
        r1=psi I(i)
        r2=psi R(j)
        r3=psi I(j)
        psi R(i)=(r0+r2)*c
        psi I(i)=(r1+r3)*c
        psi R(j)=(r0-r2)*c
        psi I(j)=(r1-r3)*c
 end subroutine
```

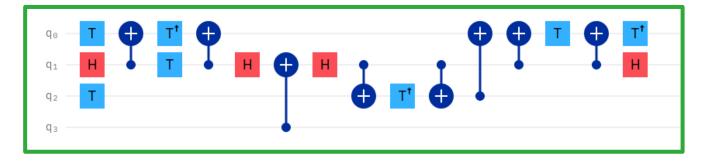
CUDA implementation



```
\begin{aligned} &\texttt{i1} = 000000000010 \cdots 0 \\ &\texttt{k} = rrrrrrrr * * * \cdots * \\ &\texttt{i} = rrrrrrrr * 0 * \cdots * \\ &\texttt{j} = rrrrrrrr * 1 * \cdots * \end{aligned}
```

```
attributes(global) subroutine H operation GPU(nstates,psi R,psi I,ie,i1,c)
        USE cudafor
        implicit real(kind=8) (a-h,o-z)
        integer(kind=8),parameter :: one=1
        integer(kind=8), value :: nstates in
        integer(kind=8) :: m1,nm1,i,j,k
        integer(kind=4), value :: i0
        real(kind=8). Yatue :: c
        real(kind=8), dimension(0:nstates-1), device:: psi R, psi I
        k=blockidx%x-1
        k=k*blockdim%x+threadidx%x-1
        m1=i1-1
        nm1=not(m1)
        i=ishft(iand(k,nm1),one)+iand(k,m1)
        j=i+i1
        r0=psi R(i)
        r1=psi I(i)
        r2=psi R(j)
        r3=psi I(j)
        psi R(i)=(r0+r2)*c
        psi I(i)=(r1+r3)*c
        psiR(j)=(r0-r2)*c
        psi I(j)=(r1-r3)*c
 end subroutine
```

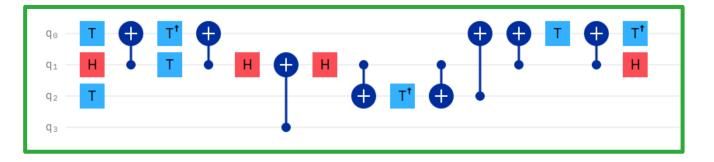
CUDA implementation



```
i1 = 000000000010 \cdots 0
k = rrrrrrrr * * * \cdots *
i = rrrrrrrr * 0 * \cdots *
j = rrrrrrrr * 1 * \cdots *
```

```
attributes(global) subroutine H operation GPU(nstates,psi R,psi I,i0,i1,c)
        USE cudafor
        implicit real(kind=8) (a-h,o-z)
        integer(kind=8),parameter :: one=1
        integer(kind=8), value :: nstates,i1
        integer(kind=8) :: m1,nm1,i,j,k
        integer(kind=4), value :: i0
        real(kind=8), value :: c
        real(kind=8),dimension(0:nstates-1),device:: psi R,psi I
        k=blockidx%x-1
        k=k*blockdim%x+threadidx%x-1
        m1=i1-1
        nm1=not(m1)
        i=ishft(iand(k,nm1),one)+iand(k,m1)
        j=i+i1
        r0=psi R(i)
        r1=psi I(i)
        r2=psi R(j)
        r3=psi I(j)
        psi R(i)=(r0+r2)*c
        psi I(i)=(r1+r3)*c
        psiR(j)=(r0-r2)*c
        psi I(j)=(r1-r3)*c
 nd subroutine
```

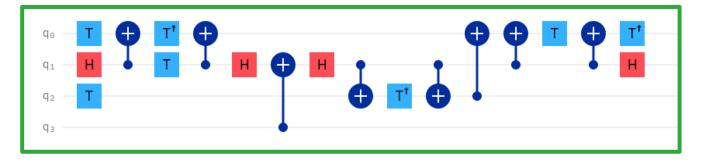
CUDA implementation



```
\begin{array}{ll} \mathtt{i1} = 000000000010\cdots 0 \\ \mathtt{k} = rrrrrrrr * * * \cdots * \\ \mathtt{i} = rrrrrrrr * 0 * \cdots * \\ \mathtt{j} = rrrrrrrr * 1 * \cdots * \end{array}
```

```
attributes(global) subroutine H operation GPU(nstates,psi R,psi I,i0,i1,c)
        USE cudafor
        implicit real(kind=8) (a-h,o-z)
        integer(kind=8),parameter :: one=1
        integer(kind=8), value :: nstates,i1
        integer(kind=8) :: m1,nm1,i,j,k
        integer(kind=4), value :: i0
        real(kind=8), value :: c
        real(kind=8),dimension(0:nstates-1),device:: psi R,psi I
        k=blockidx%x-1
        k=k*blockdim%x+threadidx%x-1
        m1=i1-1
        nm1=not(m1)
        i=ishft(iand(k,nm1),one)+iand(k,m1)
        j=i+i1
        r0=psi R(i)
        rl=psi I(i)
        r2=psi R(j)
        r3=psi I(j)
        psi R(i)=(r0+r2)*c
        psi I(i)=(r1+r3)*c
        psiR(j)=(r0-r2)*c
        psi I(j)=(r1-r3)*c
 nd subroutine
```

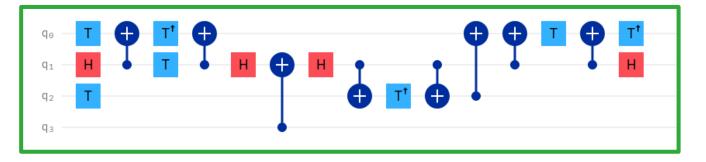
CUDA implementation



```
\begin{aligned} &\texttt{i1} = 000000000010 \cdots 0 \\ &\texttt{k} = rrrrrrrr * * * * \cdots * \\ &\texttt{i} = rrrrrrrr * 0 * \cdots * \\ &\texttt{j} = rrrrrrrr * 1 * \cdots * \end{aligned}
```

```
attributes(global) subroutine H operation GPU(nstates,psi R,psi I,i0,i1,c)
        USE cudafor
        implicit real(kind=8) (a-h,o-z)
        integer(kind=8),parameter :: one=1
        integer(kind=8), value :: nstates,i1
        integer(kind=8) :: m1,nm1,i,j,k
        integer(kind=4), value :: i0
        real(kind=8), value :: c
        real(kind=8),dimension(0:nstates-1),device:: psi R,psi I
        k=blockidx%x-1
        k=k*blockdim%x+threadidx%x-1
        m1=i1-1
        nm1=not(m1)
        i=ishft(iand(k,nm1),one)+iand(k,m1)
        j=i+i1
        r0=psi R(i)
        r1=psi I(i)
        r2=psi R(j)
        psi R(i)=(r0+r2)*c
        psi I(i)=(r1+r3)*c
        psi R(j)=(r0-r2)*c
        psi I(j)=(r1-r3)*c
 nd subroutine
```

CUDA implementation

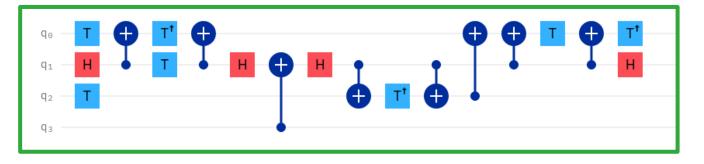


How to implement these matrix-vector multiplications in the most efficient way?

```
\psi_{00}0000000 0\cdots0
   GPU rank 0
\psi_{0000000001\cdots 1}
\psi_{00000001\,0\cdots0}
  GPU rank 1
\psi_{00000001\,1\cdots 1}
\psi_{1111111110\cdots 0}
GPU rank 255
```



CUDA implementation



How to implement these matrix-vector multiplications in the most efficient way?

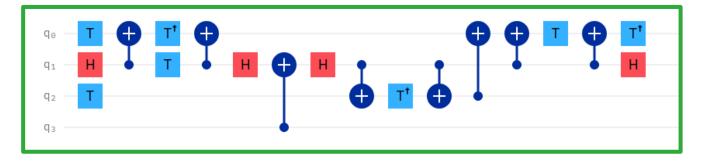
Full quantum state:

```
\overline{\psi_{00}}_{00000000000000}
   GPU rank 0
\psi_{0000000001\cdots 1}
\psi_{00000001\,0\cdots0}
   GPU rank 1
\psi_{00000001\,1\cdots 1}
\psi_{1111111110\cdots0}
GPU rank 255
```

Quantum gate on **global** qubits:



CUDA implementation



How to implement these matrix-vector multiplications in the most efficient way?

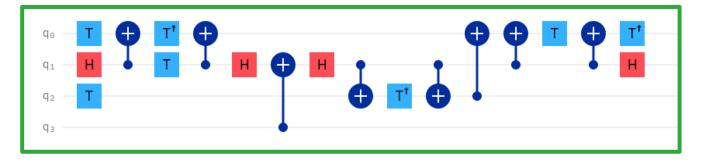
Full quantum state:

```
\overline{\psi_{00000000}} 0...0
   GPU rank 0
\psi_{0000000001\cdots 1}
\psi_{00000001\,0\cdots0}
   GPU rank 1
\psi_{00000001\,1\cdots 1}
\psi_{1111111110\cdots0}
 GPU rank 255
```

Quantum gate on **global** qubits: e.g. $\stackrel{\text{H}}{\text{H}}$ on qubit q_{32}

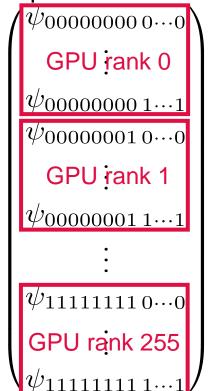


CUDA implementation



How to implement these matrix-vector multiplications in the most efficient way?

Full quantum state:



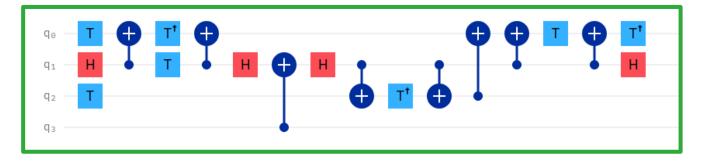
Quantum gate on **global** qubits:

e.g.
$${\color{red} \hspace{-0.07cm} \hspace{-0.0c$$

➤ Need to perform 2x2 updates of the form

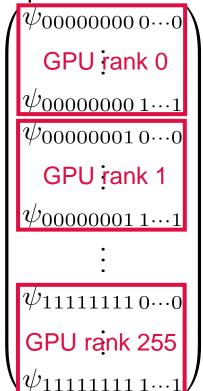
$$\begin{array}{c} \mathbf{H} & \begin{pmatrix} \psi_{********0} * \cdots * \\ \psi_{******1} * \cdots * \end{pmatrix} \end{array}$$

CUDA implementation



How to implement these matrix-vector multiplications in the most efficient way?

Full quantum state:



Quantum gate on **global** qubits:

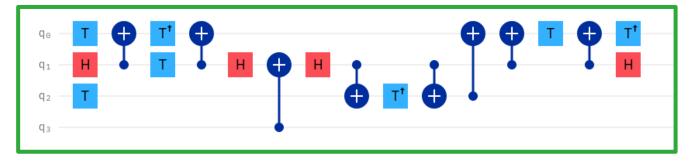
e.g.
$$f H$$
 on qubit q_{32}

➤ Need to perform 2x2 updates of the form

$$\begin{array}{c} \mathbf{H} & \begin{pmatrix} \psi_{*******0} * \cdots * \\ \psi_{******1} * \cdots * \end{pmatrix} \end{array}$$

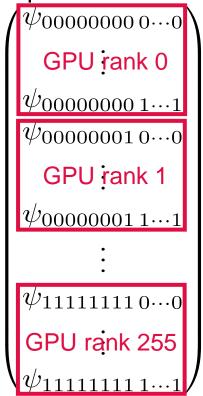
> Problem: the numbers are on separate GPUs

CUDA implementation



How to implement these matrix-vector multiplications in the most efficient way?

Full quantum state:



Quantum gate on **global** qubits:

e.g.
$${\color{red} \hspace{-0.07cm} \hspace{-0.0c$$

➤ Need to perform 2x2 updates of the form

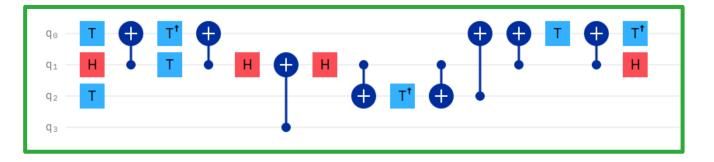
$$\begin{array}{c} \mathbf{H} & \begin{pmatrix} \psi_{*******0} * \cdots * \\ \psi_{******1} * \cdots * \end{pmatrix} \end{array}$$

Problem: the numbers are on separate GPUs

Page 12

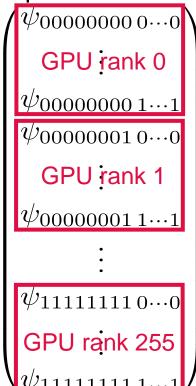
- Naïve solution:
 - > Transfer $2^n/2 \psi's$ (8 TiB), perform \blacksquare , transfer back

CUDA implementation



How to implement these matrix-vector multiplications in the most efficient way?

Full quantum state:



Quantum gate on **global** qubits:

e.g.
$$\blacksquare$$
 on qubit q_{32}

➤ Need to perform 2x2 updates of the form

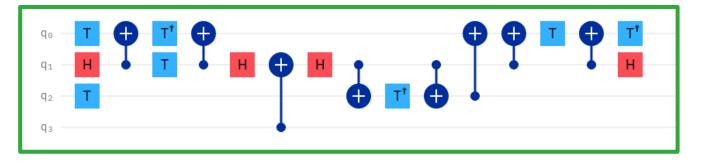
$$\begin{array}{c} \mathbf{H} & \begin{pmatrix} \psi_{*******0} * \cdots * \\ \psi_{******1} * \cdots * \end{pmatrix} \end{array}$$

Problem: the numbers are on separate GPUs

Page 12

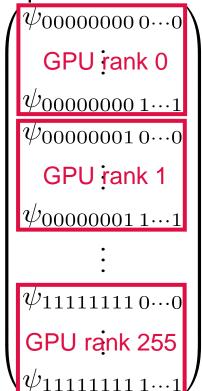
- Naïve solution:
 - ightharpoonup Transfer $2^n/2~\psi'\mathrm{s}$ (8 TiB), perform \blacksquare , transfer back
- > Optimal solution:
 - \triangleright Exchange global and local qubit, e.g. $q_{32} \leftrightarrow q_0$

CUDA implementation



How to implement these matrix-vector multiplications in the most efficient way?

Full quantum state:



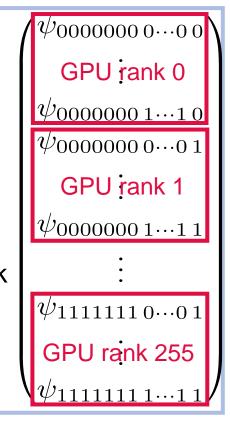
Quantum gate on **global** qubits:

e.g. \blacksquare on qubit q_{32}

> Need to perform 2x2 updates of the form

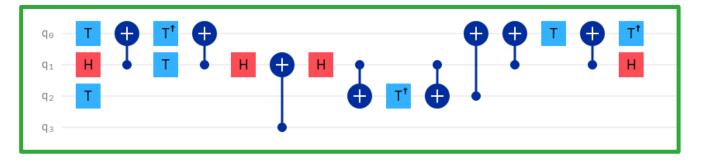
$$\begin{array}{c} \mathbf{H} & \begin{pmatrix} \psi_{*******0} * \cdots * \\ \psi_{******1} * \cdots * \end{pmatrix} \end{array}$$

- > Problem: the numbers are on separate GPUs
- Naïve solution:
 - > Transfer $2^n/2 \psi's$ (8 TiB), perform \blacksquare , transfer back
- Optimal solution:
 - \triangleright Exchange global and local qubit, e.g. $q_{32} \leftrightarrow q_0$



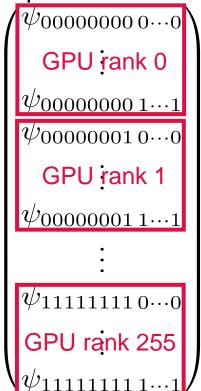


CUDA implementation



How to implement these matrix-vector multiplications in the most efficient way?

Full quantum state:



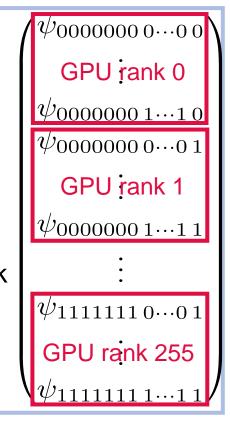
Quantum gate on **global** qubits:

e.g.
$$\blacksquare$$
 on qubit q_{32}

> Need to perform 2x2 updates of the form

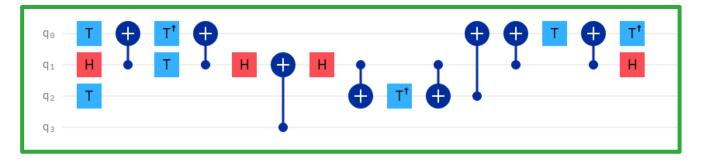
$$\begin{array}{c} \mathbf{H} \begin{pmatrix} \psi_{rrrrrr0*\cdots*r} \\ \psi_{rrrrr1*\cdots*r} \end{pmatrix} \end{array}$$

- > Problem: the numbers are on separate GPUs
- Naïve solution:
 - > Transfer $2^n/2 \psi's$ (8 TiB), perform \blacksquare , transfer back
- Optimal solution:
 - \triangleright Exchange global and local qubit, e.g. $q_{32} \leftrightarrow q_0$



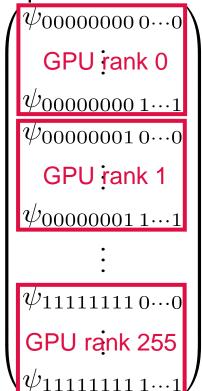


CUDA implementation



How to implement these matrix-vector multiplications in the most efficient way?

Full quantum state:



Quantum gate on **global** qubits:

e.g. \blacksquare on qubit q_{32}

> Need to perform 2x2 updates of the form

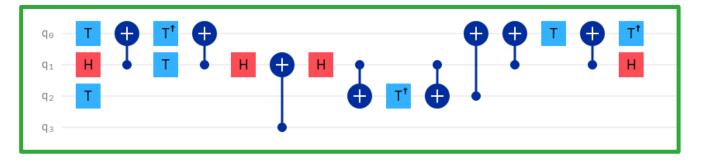
$$\begin{array}{c} \mathbf{H} \begin{pmatrix} \psi_{rrrrrr0*\cdots*r} \\ \psi_{rrrrr1*\cdots*r} \end{pmatrix} \end{array}$$

- > Problem: the numbers are on separate GPUs
- Naïve solution:
 - > Transfer $2^n/2 \ \psi's$ (8 TiB), perform \blacksquare , transfer back
- Optimal solution:
 - \triangleright Exchange global and local qubit, e.g. $q_{32} \leftrightarrow q_0$
 - ➤ Keep track of qubit assignment in a permutation

```
\psi_{0000000000\cdots00}
   GPU rank 0
\psi_{000000001\cdots 10}
\psi_{0000000000\cdots01}
   GPU rank 1
\psi_{000000001\cdots11}
\psi_{111111110\cdots01}
 GPU rank 255
\psi_{111111111\cdots 1} 1
```

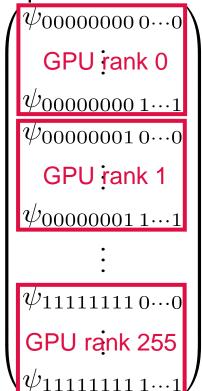


CUDA implementation



How to implement these matrix-vector multiplications in the most efficient way?

Full quantum state:



Quantum gate on **global** qubits:

e.g. \mathbf{H} on qubit q_{32}

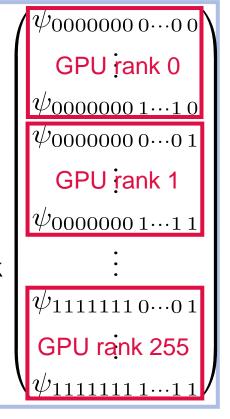
Need to perform 2x2 updates of the form

$$\begin{array}{c} \mathbf{H} \begin{pmatrix} \psi_{rrrrrr0*\cdots*r} \\ \psi_{rrrrr1*\cdots*r} \end{pmatrix} \end{array}$$

> Problem: the numbers are on separate GPUs

Page 12

- Naïve solution:
 - > Transfer $2^n/2 \ \psi's$ (8 TiB), perform \blacksquare , transfer back
- > Optimal solution:
 - \triangleright Exchange global and local qubit, e.g. $q_{32} \leftrightarrow q_0$
 - Keep track of qubit assignment in a permutation
 - > Transfer $2^n/2 \psi'$ s only **once**



MPI communication scheme

Before transfer:

```
\begin{pmatrix} \psi_{rrrrrrrr*\cdots*} \\ \psi_{rrrrrrr*\cdots*} \end{pmatrix}
```

```
call psi2MPIbuf<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_I)
istat=cudaStreamSynchronize()
call MPI_SENDRECV(DD(0)%buf_I,m0,MPI_REAL8,idest,0,DD(0)%buf_R,m0,MPI_REAL8,idest,MPI_ANY_TAG,MPI_COMM_WORLD,MPI_STATUS_IGNORE,IERR)
call MPIbuf2psi<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_R)
istat=cudaStreamSynchronize()
```

MPI communication scheme

Before transfer:

$$\begin{pmatrix} \psi_{rrrrrrr*\cdots*} \\ \psi_{rrrrrrr*\cdots*} \end{pmatrix}$$

After transfer:

```
\begin{pmatrix} \psi_{rrrrrr*\cdots*r} \\ \psi_{rrrrrr*\cdots*r} \end{pmatrix}
```

```
call psi2MPIbuf<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_I)
istat=cudaStreamSynchronize()
call MPI_SENDRECV(DD(0)%buf_I,m0,MPI_REAL8,idest,0,DD(0)%buf_R,m0,MPI_REAL8,idest,MPI_ANY_TAG,MPI_COMM_WORLD,MPI_STATUS_IGNORE,IERR)
call MPIbuf2psi<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_R)
istat=cudaStreamSynchronize()
```

MPI communication scheme

Before transfer:

$$\begin{pmatrix} \psi_{rrrrrrrr*\cdots*} \\ \psi_{rrrrrrr*\cdots*} \end{pmatrix}$$

After transfer:

$$\begin{pmatrix} \psi_{rrrrrrr*\cdots*r} \\ \psi_{rrrrrr*\cdots*r} \end{pmatrix}$$

```
call psi2MPIbuf<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_I)
istat=cudaStreamSynchronize()
call MPI_SENDRECV(DD(0)%buf_I,m0,MPI_REAL8,idest,0,DD(0)%buf_R,m0,MPI_REAL8,idest,MPI_ANY_TAG,MPI_COMM_WORLD,MPI_STATUS_IGNORE,IERR)
call MPIbuf2psi<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_R)
istat=cudaStreamSynchronize()
```

MPI communication scheme

Before transfer:

$$\left(egin{array}{c} \psi_{rrrrrrr*\cdots*} \ \psi_{rrrrrrr*\cdots*} \end{array}
ight)$$

After transfer:

$$\begin{pmatrix} \psi_{rrrrrrr*\cdots*r} \\ \psi_{rrrrrr*\cdots*r} \end{pmatrix}$$

```
\psi_{000000000000...00}
\psi_{00000000000\cdots01}
\psi_{00000000000\cdots 10}
\psi_{00000000000\cdots11}
    GPU rank 0
\psi_{000000000} 1\cdots11
\psi_{00000001\,0\cdots00}
\psi_{00000001\,0\cdots01}
\psi_{000000010\cdots 10}
\psi_{00000001\,0\cdots11}
   GPU rank 1
\psi_{000000001\,1\cdots11}
```

```
call psi2MPIbuf<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_I)
istat=cudaStreamSynchronize()
call MPI_SENDRECV(DD(0)%buf_I,m0,MPI_REAL8,idest,0,DD(0)%buf_R,m0,MPI_REAL8,idest,MPI_ANY_TAG,MPI_COMM_WORLD,MPI_STATUS_IGNORE,IERR)
call MPIbuf2psi<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_R)
istat=cudaStreamSynchronize()
```

MPI communication scheme

Before transfer:

$$\begin{pmatrix} \psi_{rrrrrrr*\cdots*} \\ \psi_{rrrrrrr*\cdots*} \end{pmatrix}$$

After transfer:

$$\begin{pmatrix} \psi_{rrrrrrr*\cdots*r} \ \psi_{rrrrrrr*\cdots*r} \end{pmatrix}$$

```
\psi_{00000000000\cdots01}
                             GPU0
\psi_{000000000000...10}
                              buffer
\psi_{00000000000\cdots11}
\psi_{0000000001\cdots11}
\psi_{00000001\,0\cdots00}
\psi_{00000001\,0\cdots01}
\psi_{000000010\cdots 10}
\psi_{00000001\,0\cdots11}
   GPU rank 1
\psi_{000000001\,1\cdots11}
```

```
call psi2MPIbuf<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_I)
istat=cudaStreamSynchronize()
call MPI_SENDRECV(DD(0)%buf_I,m0,MPI_REAL8,idest,0,DD(0)%buf_R,m0,MPI_REAL8,idest,MPI_ANY_TAG,MPI_COMM_WORLD,MPI_STATUS_IGNORE,IERR)
call MPIbuf2psi<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_R)
istat=cudaStreamSynchronize()
```

MPI communication scheme

Before transfer:

$$\begin{pmatrix} \psi_{rrrrrrrr*\cdots*} \\ \psi_{rrrrrrr*\cdots*} \end{pmatrix}$$

After transfer:

$$\left(egin{array}{c} \psi_{rrrrrrr*\cdots*r} \ \psi_{rrrrrr*\cdots*r} \end{array}
ight)$$

```
\psi_{000000000000...01}
                         GPU0
buffer
\psi_{00000000000\cdots11}
\psi_{0000000001\cdots11}
\psi_{00000001\,0\cdots00}
\psi_{000000010...01}
                         GPU1
\psi_{000000010\cdots 10}
                         buffer
\psi_{000000010\cdots 11}
\psi_{00}0000011...11
```

```
call psi2MPIbuf<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_I)
istat=cudaStreamSynchronize()
call MPI_SENDRECV(DD(0)%buf_I,m0,MPI_REAL8,idest,0,DD(0)%buf_R,m0,MPI_REAL8,idest,MPI_ANY_TAG,MPI_COMM_WORLD,MPI_STATUS_IGNORE,IERR)
call MPIbuf2psi<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_R)
istat=cudaStreamSynchronize()
```

MPI communication scheme

Before transfer:

$$\begin{pmatrix} \psi_{rrrrrrr*\cdots*} \\ \psi_{rrrrrrr*\cdots*} \end{pmatrix}$$

After transfer:

$$\left(\psi_{rrrrrr*\cdots*r}
ight) \ \left(\psi_{rrrrrr*\cdots*r}
ight)$$

- > Build buffer on each GPU
- > Transfer with MPI

```
\psi_{000000000000...01}
                            GPU0
\psi_{000000000000...10}
                            buffer
\psi_{00000000000\cdots11}
\psi_{0000000001\cdots11}
\psi_{000000001\,0\cdots00}
\psi_{000000010...01}
                            GPU1
\psi_{000000010\cdots 10}
                            buffer
\psi_{000000010...11}
\psi_{00000001} 1...11
```

```
call psi2MPIbuf<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_I)
istat=cudaStreamSynchronize()
call MPI_SENDRECV(DD(0)%buf_I,m0,MPI_REAL8,idest,0,DD(0)%buf_R,m0,MPI_REAL8,idest,MPI_ANY_TAG,MPI_COMM_WORLD,MPI_STATUS_IGNORE,IERR)
call MPIDuT2ps1<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_R)
istat=cudaStreamSynchronize()
```

MPI communication scheme

Before transfer:

$$\begin{pmatrix} \psi_{rrrrrrr*\cdots*} \\ \psi_{rrrrrrr*\cdots*} \end{pmatrix}$$

After transfer:

$$\begin{pmatrix} \psi_{rrrrrrr*\cdots*r} \\ \psi_{rrrrrr*\cdots*r} \end{pmatrix}$$

- > Build buffer on each GPU
- > Transfer with MPI

```
\psi_{00000000000000000}
GPU0
\psi_{000000000000...10}
                          buffer
\psi_{000000000000...11}
\psi_{0000000001...11}
\psi_{000000010...00}
\psi_{000000010...01}
                          GPU1
\psi_{000000010...10}
                          buffer
\psi_{000000010...11}
\psi_{000000011...11}
```

```
n=blockidx%x-1
   n=n*blockdim%x+threadidx%x-1
   idx=n+noff
   do i=0, nswap - 1
    j=bitmasklist(i)
   idx=ishft(iand(idx,not(j)),one)+iand(idx,j)
   enddo
   buf(n)=phi(ior(idx,jsourmask))
                                       \psi_{0000000000\cdots00} 0
|\psi_{0000000000000000}|
                                          GPU rank 0
  GPU rank 0
                                       \psi_{000000001\cdots 10}
\psi_{0000000001\cdots 1}
                                       |\psi_{00000000000...01}|
|\psi_{00000001}\,{}_{0\cdots0}|
  GPU rank 1
                                          GPU rank 1
\psi_{00000001\,1\cdots 1}
                                       \psi_{000000001...11}
                                       |\psi_{11111111\,0\cdots 0\,1}|
\psi_{1111111110\cdots0}
GPU rank 255
                                        GPU rank 255
\psi111111111\cdots1
                                      \Psi_{1111111111\cdots 1}
```

```
call psi2MPIbuf<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_I)
istat=cudaStreamSynchronize()
call MPI_SENDRECV(DD(0)%buf_I,m0,MPI_REAL8,idest,0,DD(0)%buf_R,m0,MPI_REAL8,idest,MPI_ANY_TAG,MPI_COMM_WORLD,MPI_STATUS_IGNORE,IERR)
call MPIDuT2ps1<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_R)
istat=cudaStreamSynchronize()
```

MPI communication scheme

Before transfer:

$$\begin{pmatrix} \psi_{rrrrrrr*\cdots*} \\ \psi_{rrrrrrr*\cdots*} \end{pmatrix}$$

After transfer:

$$\begin{pmatrix} \psi_{rrrrrrr*\cdots*r} \\ \psi_{rrrrrr*\cdots*r} \end{pmatrix}$$

- ➤ Build buffer on each GPU
- Transfer with MPI
- > Retrieve from buffer

```
\psi_{000000000000...10}
\psi_{000000000000...11}
\psi_{0000000001...11}
\psi_{000000010...00}
\psi_{000000010...01}
\psi_{000000010...10}
\psi_{000000010...11}
\psi_{000000011...11}
```

```
n=blockidx%x-1
   n=n*blockdim%x+threadidx%x-1
   idx=n+noff
   do i=0, nswap - 1
    j=bitmasklist(i)
   idx=ishft(iand(idx,not(j)),one)+iand(idx,j)
   enddo
   buf(n)=phi(ior(idx,jsourmask))
                                      \psi_{0000000000\cdots00} 0
\psi_{000000000000000}
  GPU rank 0
                                         GPU rank 0
                                      \psi_{000000001\cdots 10}
\psi_{0000000001\cdots 1}
                                       |\psi_{00000000000...01}|
|\psi_{00000001}\,{}_{0\cdots0}|
  GPU rank 1
                                         GPU rank 1
\psi_{00000001\,1\cdots1}
                                       \psi_{000000001...11}
                                      \psi_{111111110\cdots01}
\psi_{1111111110\cdots0}
GPU rank 255
                                       GPU rank 255
\psi11111111_{1\cdots 1}
                                     \Psi_{1111111111\cdots 1}
```

```
call psi2MPIbuf<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_I)
istat=cudaStreamSynchronize()
call MPI SENDRECV(DD(0)%buf_I,m0,MPI_REAL8,idest,0,DD(0)%buf_R,m0,MPI_REAL8,idest,MPI_ANY_TAG,MPI_COMM_WORLD,MPI_STATUS_IGNORE,IERR)
call MPIbuf2psi<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_R)
istat=cudaStreamSynchronize()
```

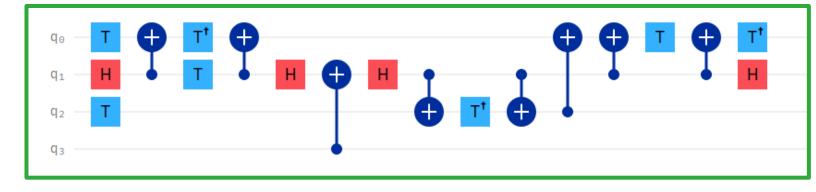
GPU0

buffer

GPU1

buffer

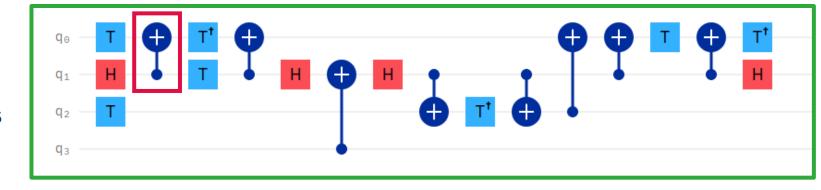
MPI communication: Two-qubit gates





MPI communication: Two-qubit gates

➤ CNOT gate (controlled-NOT):





MPI communication: Two-qubit gates

- q0
 T
 T
 T
 T
 T
 T
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
 H
- ➤ CNOT gate (controlled-NOT):
 - > Swap all coefficients where control qubit C is 1

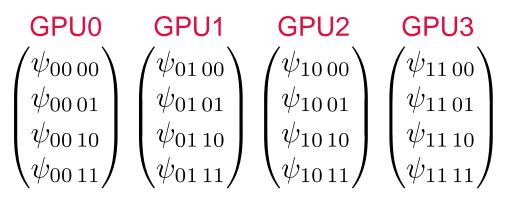
С	${ m T}$	CT	CT	CT
0	0	01	10	11
/ -	1	0	0	0 \
)	1	0	0
)	0	0	1
/ ()	0	1	0 /

MPI communication: Two-qubit gates

- ➤ CNOT gate (controlled-NOT):
 - > Swap all coefficients where control qubit C is 1

CT	CT	CT	CT
00	01	10	11
/ 1	0	0	0 \
0	1	0	0
0	0	0	1
0	0	1	0 /
	00 / 1 0	$\begin{array}{cccc} 00 & 01 \\ 1 & 0 \\ 0 & 1 \\ \end{array}$	$\begin{array}{cccccccccccccccccccccccccccccccccccc$

$$\psi_{q_3q_2q_1q_0}$$



MPI communication: Two-qubit gates

- ➤ CNOT gate (controlled-NOT):
 - > Swap all coefficients where control qubit C is 1

CT	CT	CT	CT
00	01	10	11
/ 1	0	0	0 \
0	1	0	0
0	0	0	1
$\int 0$	0	1	$_{0}$ /

 $\psi_{q_3q_2q_1q_0}$ \dagger \dagger C T

 $\begin{pmatrix} \psi_{00\ 00} \\ \psi_{00\ 01} \\ \psi_{00\ 10} \\ \psi_{00\ 11} \end{pmatrix} \begin{pmatrix} \psi_{01\ 00} \\ \psi_{01\ 01} \\ \psi_{01\ 10} \\ \psi_{01\ 11} \end{pmatrix}$

 $\begin{pmatrix} \psi_{01\,00} \\ \psi_{01\,01} \\ \psi_{01\,10} \\ \psi_{01\,11} \end{pmatrix} \begin{pmatrix} \psi_{10\,00} \\ \psi_{10\,01} \\ \psi_{10\,10} \\ \psi_{10\,11} \end{pmatrix}$

GPU2

 $\begin{pmatrix} \psi_{11\,00} \\ \psi_{11\,01} \\ \psi_{11\,10} \\ \psi_{11\,11} \end{pmatrix}$

GPU3

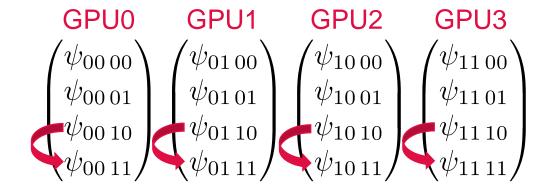
➤ "2 local": GPUs swap locally

MPI communication: Two-qubit gates

- ➤ CNOT gate (controlled-NOT):
 - > Swap all coefficients where control qubit C is 1

CT	CT	CT	CT
00	01	10	11
/ 1	0	0	0 \
0	1	0	0
0	0	0	1
$\int 0$	0	1	0 /

 $\psi_{q_3q_2q_1q_0}$ $\uparrow \uparrow$ $\mathsf{C}\,\mathsf{T}$



➤ "2 local": GPUs swap locally

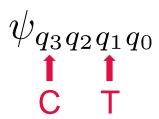


MPI communication: Two-qubit gates

 q_2 q_3

- ➤ CNOT gate (controlled-NOT):
 - > Swap all coefficients where control qubit C is 1

CT	CT	CT	CT
00	01	10	11
/ 1	0	0	0 \
0	1	0	0
0	0	0	1
$\int 0$	0	1	0 /



Page 14

 $\psi_{00\,00}$ $\psi_{00\,01}$ $\psi_{00\,10}$

GPU0

 $\psi_{01\,00}$ $\psi_{01\,01}$ $\psi_{01\,10}$

GPU1

GPU2 GPU3 $\psi_{10\,00}$ $\psi_{11\,00}$ $\psi_{10\,01}$ $\psi_{11\,01}$ $\psi_{10\,10}$ $\psi_{11\,10}$

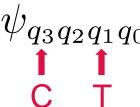
- ➤ "2 local": GPUs swap locally
- ➤ "C global, T local": GPUs swap locally

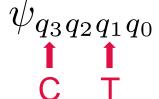
MPI communication: Two-qubit gates

 q_2 q_3

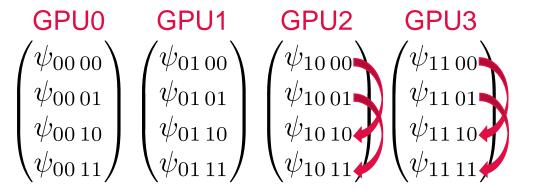
- ➤ CNOT gate (controlled-NOT):
 - > Swap all coefficients where control qubit C is 1

CT	CT	CT	CT
00	01	10	11
/ 1	0	0	0 \
0	1	0	0
0	0	0	1
/ 0	0	1	0 /





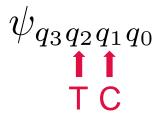
- ➤ "2 local": GPUs swap locally
- ➤ "C global, T local": GPUs swap locally



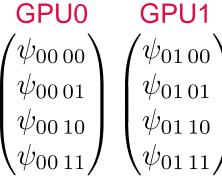
MPI communication: Two-qubit gates

- ➤ CNOT gate (controlled-NOT):
 - > Swap all coefficients where control qubit C is 1

CT	CT	CT	CT
00	01	10	11
/ 1	0	0	0 \
0	1	0	0
0	0	0	1
/0	0	1	0 /







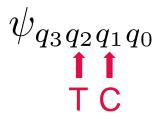
$$\begin{pmatrix} \psi_{10\,00} \\ \psi_{10\,01} \\ \psi_{10\,10} \\ \psi_{10\,11} \end{pmatrix} \begin{pmatrix} \psi_{11\,00} \\ \psi_{11\,01} \\ \psi_{11\,11} \\ \psi_{11\,11} \end{pmatrix}$$

- ➤ "2 local": GPUs swap locally
- ➤ "C global, T local": GPUs swap locally
- > "C local, T global": MPI transfer ½ of all coefficients

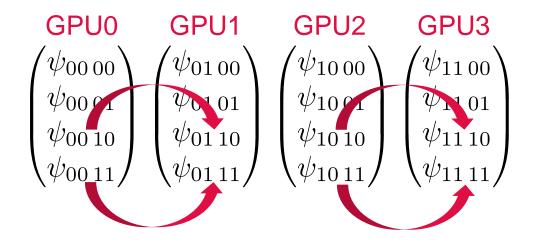
MPI communication: Two-qubit gates

- ➤ CNOT gate (controlled-NOT):
 - > Swap all coefficients where control qubit C is 1

CT	CT	CT	CT
00	01	10	11
/ 1	0	0	0 \
0	1	0	0
0	0	0	1
/0	0	1	0 /



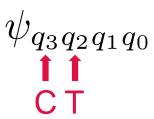
- ➤ "2 local": GPUs swap locally
- ➤ "C global, T local": GPUs swap locally
- > "C local, T global": MPI transfer ½ of all coefficients

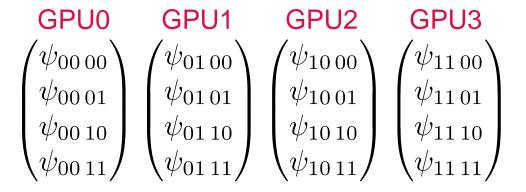


MPI communication: Two-qubit gates

- ➤ CNOT gate (controlled-NOT):
 - > Swap all coefficients where control qubit C is 1

CT	CT	CT	CT
00	01	10	11
/ 1	0	0	$0 \setminus$
0	1	0	0
0	0	0	1
/ 0	0	1	0 /



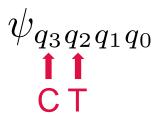


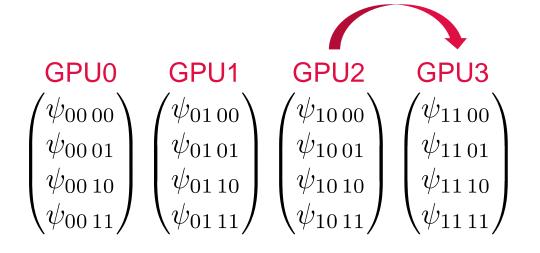
- ➤ "2 local": GPUs swap locally
- > "C global, T local": GPUs swap locally
- > "C local, T global": MPI transfer ½ of all coefficients
- ➤ "2 global": MPI transfer ½ of all coefficients (or relabel GPUs)

MPI communication: Two-qubit gates

- ➤ CNOT gate (controlled-NOT):
 - > Swap all coefficients where control qubit C is 1

CT	CT	CT	CT
00	01	10	11
/ 1	0	0	0 \
0	1	0	0
0	0	0	1
/0	0	1	0 /





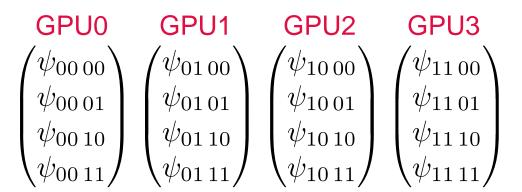
- ➤ "2 local": GPUs swap locally
- ➤ "C global, T local": GPUs swap locally
- > "C local, T global": MPI transfer ½ of all coefficients
- ➤ "2 global": MPI transfer ½ of all coefficients (or relabel GPUs)

MPI communication: Two-qubit gates

- ➤ CNOT gate (controlled-NOT):
 - > Swap all coefficients where control qubit C is 1

CT	CT	CT	CT
00	01	10	11
/ 1	0	0	0 \
0	1	0	0
0	0	0	1
/ 0	0	1	0 /

$$\psi_{q_3q_2q_1q_0}$$



- ➤ "2 local": GPUs swap locally
- > "C global, T local": GPUs swap locally
- > "C local, T global": MPI transfer ½ of all coefficients
- ➤ "2 global": MPI transfer ½ of all coefficients (or relabel GPUs)
- > For benchmarking: MPI transfer whenever coefficients on different GPUs are combined

Running on JUWELS Booster





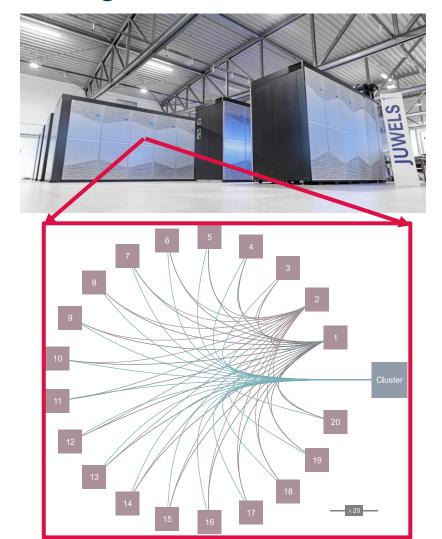
Running on JUWELS Booster



➤ 936 Compute nodes

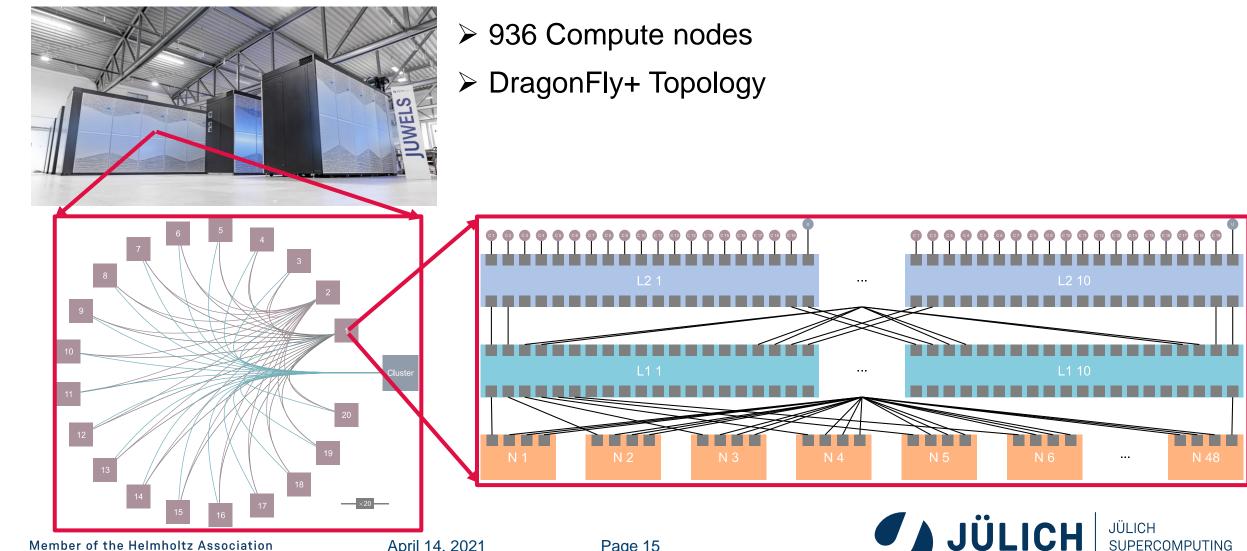


Running on JUWELS Booster



- ➤ 936 Compute nodes
- DragonFly+ Topology

Running on JUWELS Booster



CENTRE

Forschungszentrum

Running on JUWELS Booster



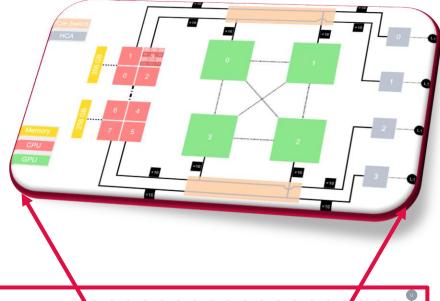
- ➤ 936 Compute nodes
- DragonFly+ Topology

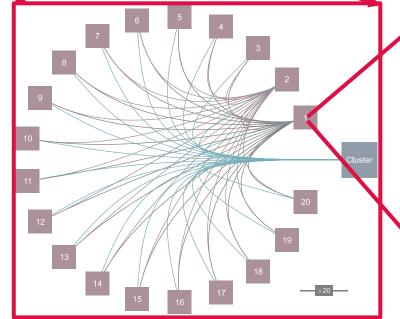
4 A100 GPUs per node

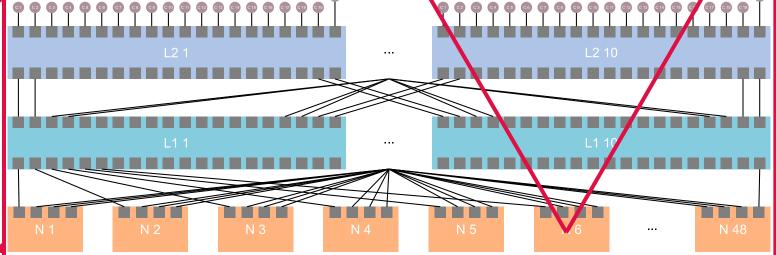
JÜLICH

CENTRE

SUPERCOMPUTING







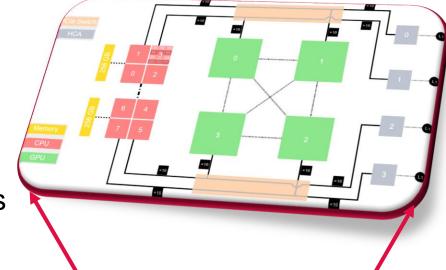
April 14, 2021

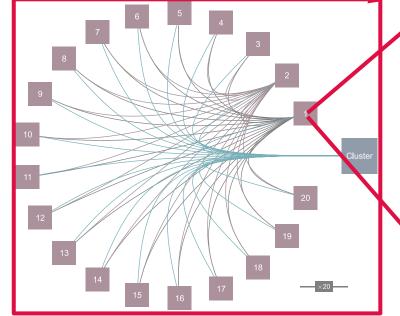
Running on JUWELS Booster

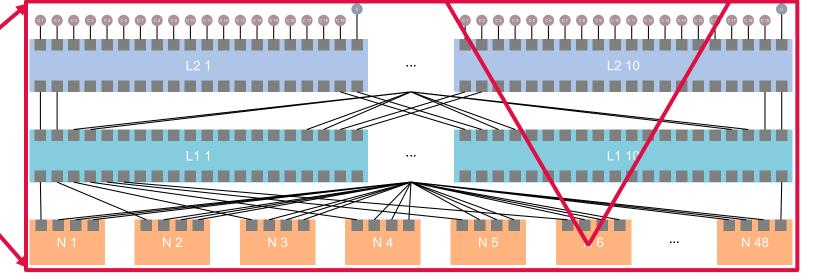


- ➤ 936 Compute nodes
- DragonFly+ Topology
- > 3744 NVIDIA A100 GPUs





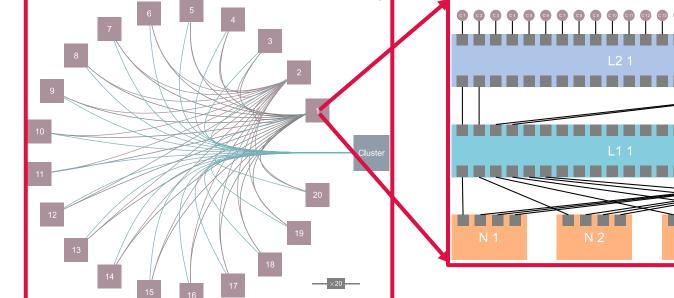


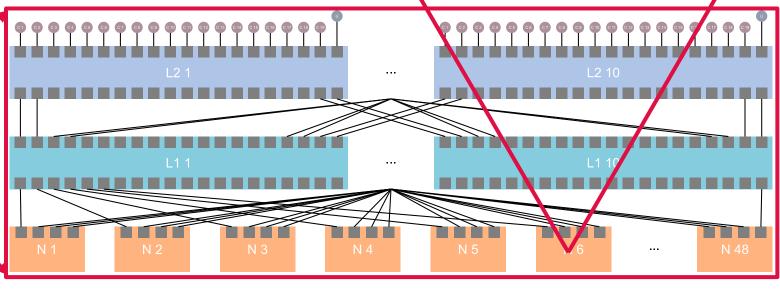


Running on JUWELS Booster



- ➤ 936 Compute nodes
- DragonFly+ Topology
- > 3744 NVIDIA A100 GPUs
- > 73 PFLOP/s





4 A100 GPUs per node

Running on JUWELS Booster

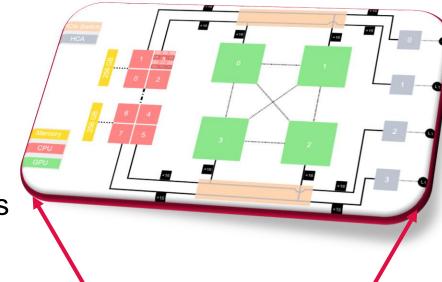


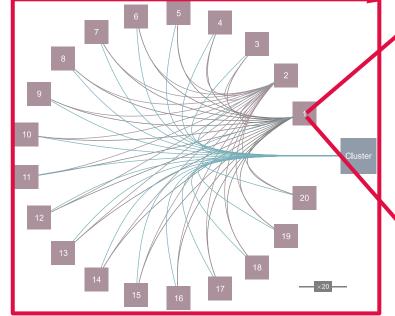
- Top500 Nov-2020:
- > #1 Europe
- > #7 World
- #3 Green500

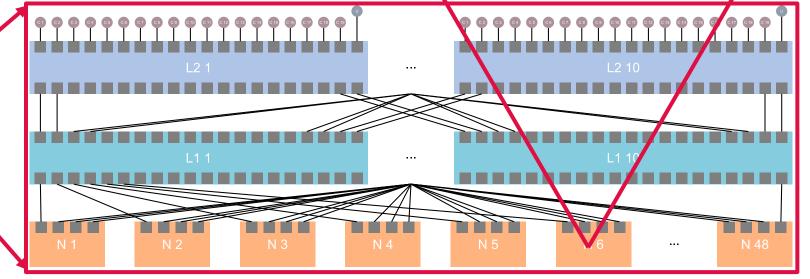


- ➤ 936 Compute nodes
- DragonFly+ Topology
- > 3744 NVIDIA A100 GPUs
- > 73 PFLOP/s

4 A100 GPUs per node







Why we can use it to benchmark GPU clusters like JUWELS Booster



Why we can use it to benchmark GPU clusters like JUWELS Booster

> Memory-intensive:

ightharpoonup For 40 qubits: $2^{40}\,\psi'\mathrm{s}$ = 16 TiB memory



Why we can use it to benchmark GPU clusters like JUWELS Booster

➤ Memory-intensive:

ightharpoonup For 40 qubits: $2^{40} \, \psi' \mathrm{s} =$ 16 TiB memory

> Network-intensive:

- > Global gates need transfer of **one half** of all memory
- > For 40 qubits: $2^{40}/2 \psi' s = 8$ TiB transfer



Why we can use it to benchmark GPU clusters like JUWELS Booster

➤ Memory-intensive:

ightharpoonup For 40 qubits: $2^{40} \, \psi' \mathrm{s} =$ 16 TiB memory

> Network-intensive:

> Global gates need transfer of **one half** of all memory

> For 40 qubits: $2^{40}/2 \, \psi' s = 8$ TiB transfer

➤ High GPU utilization

- > For 40 qubits:
 - > 32 GiB on 512 GPUs
 - > 16 GiB on 1024 GPUs
 - ➤ 8 GiB on 2048 GPUs



Why we can use it to benchmark GPU clusters like JUWELS Booster

➤ Memory-intensive:

> For 40 qubits: $2^{40} \psi' s = 16$ TiB memory

> Network-intensive:

- > Global gates need transfer of **one half** of all memory
- > For 40 qubits: $2^{40}/2 \psi' s = 8$ TiB transfer

> High GPU utilization

- > For 40 qubits:
 - > 32 GiB on 512 GPUs
 - > 16 GiB on 1024 GPUs
 - > 8 GiB on 2048 GPUs

QPU = Quantum Processing Unit

➤ Using **GPUs** to simulate **universal QPUs**



Why we can use it to benchmark GPU clusters like JUWELS Booster

➤ Memory-intensive:

> For 40 qubits: $2^{40} \psi' s = 16$ TiB memory

> Network-intensive:

- > Global gates need transfer of **one half** of all memory
- > For 40 qubits: $2^{40}/2 \psi' s = 8$ TiB transfer

≻ High GPU utilization

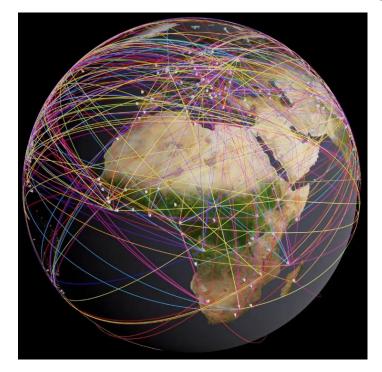
- > For 40 qubits:
 - > 32 GiB on 512 GPUs
 - > 16 GiB on 1024 GPUs
 - > 8 GiB on 2048 GPUs

QPU = Quantum Processing Unit

➤ Using GPUs to simulate universal QPUs

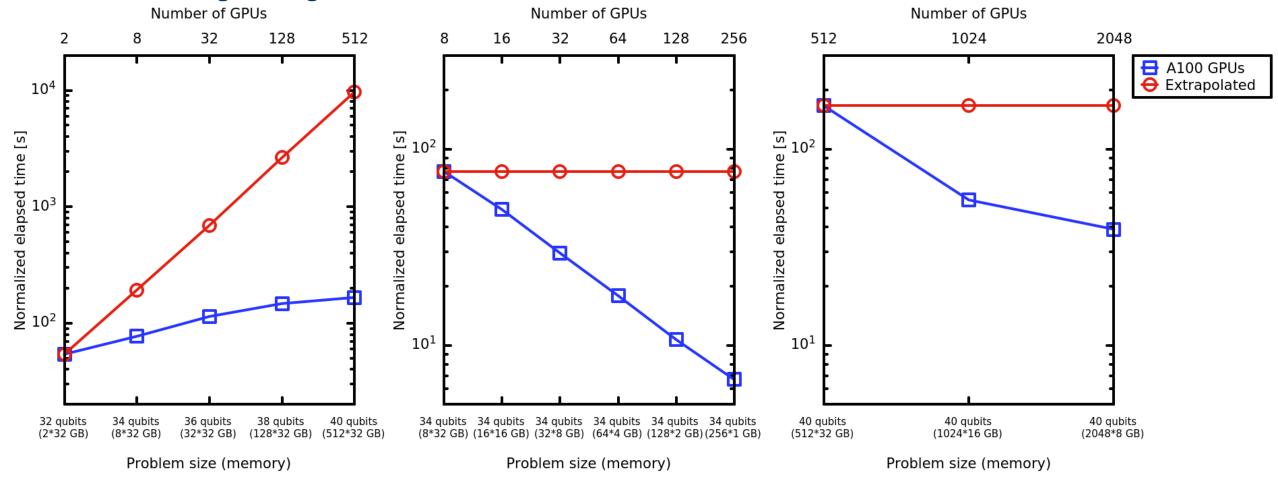
Application:

Solve airplane scheduling problems using the Quantum Approximate Optimization Algorithm



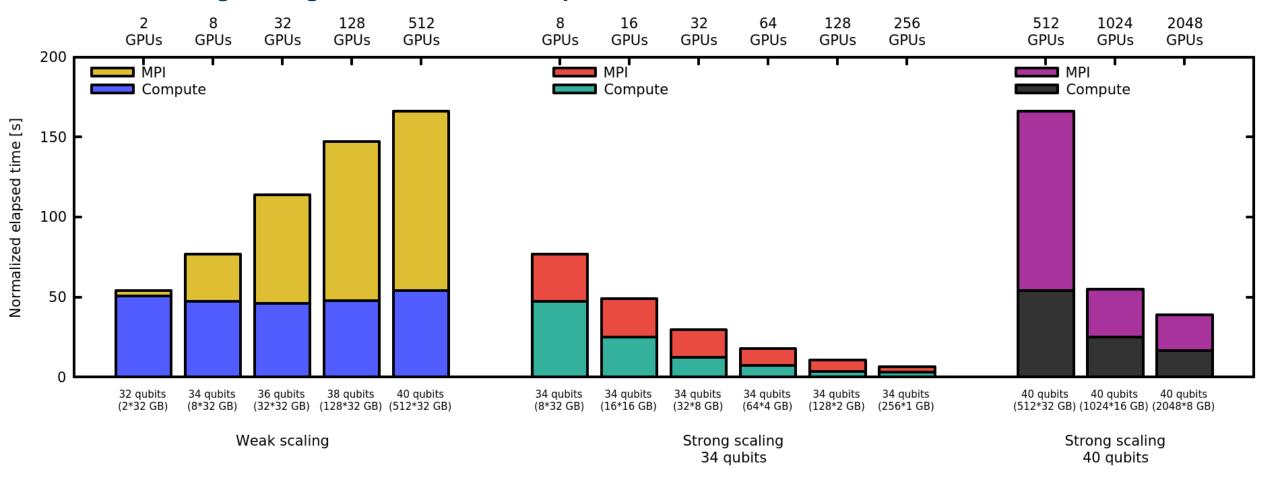


Weak and strong scaling results

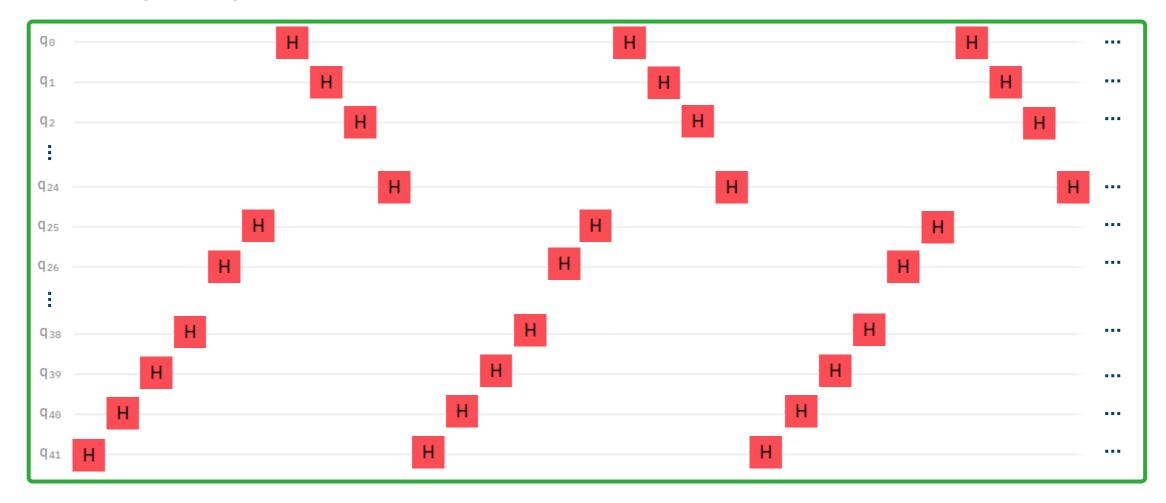




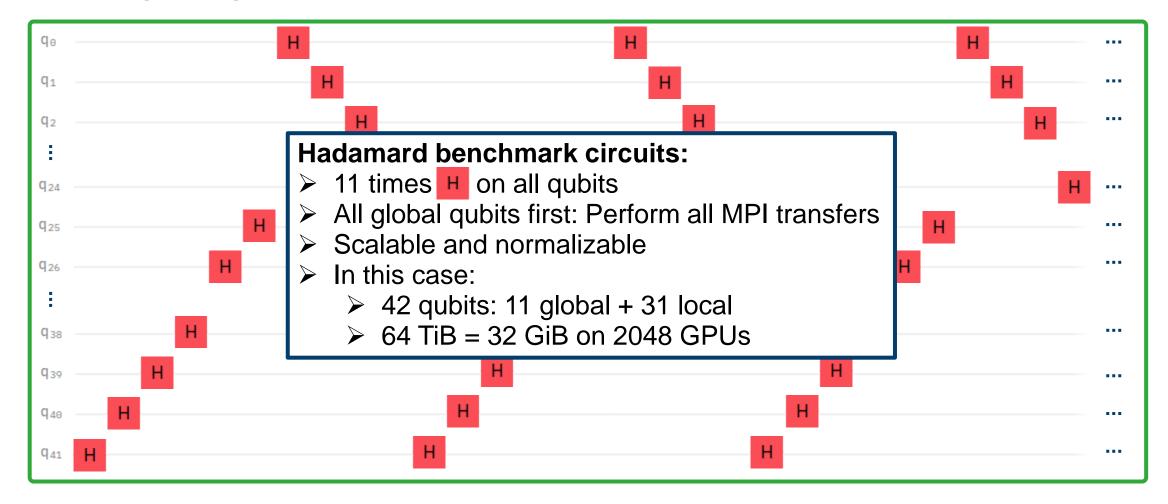
Weak and strong scaling results: MPI vs. Compute Time

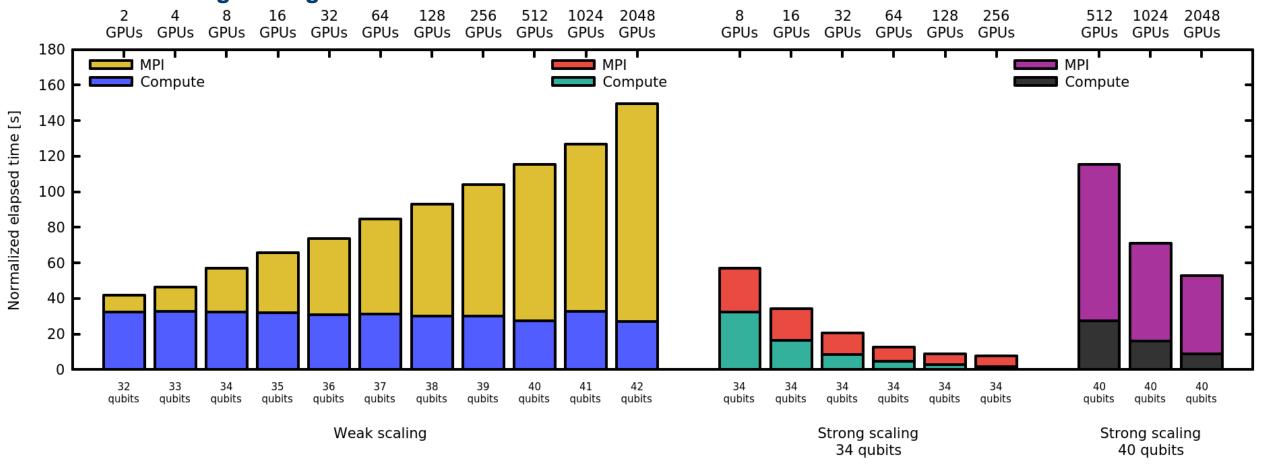




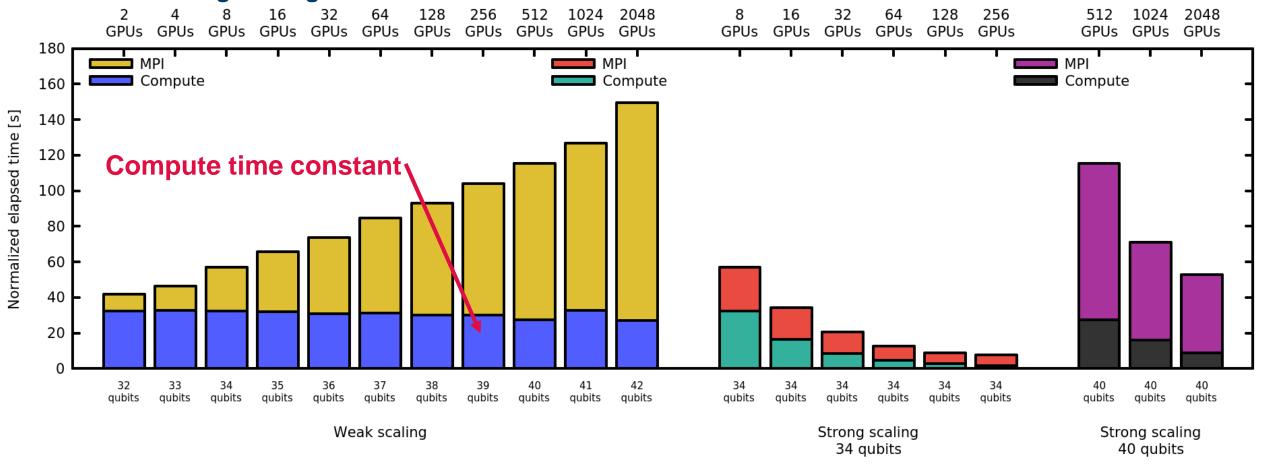


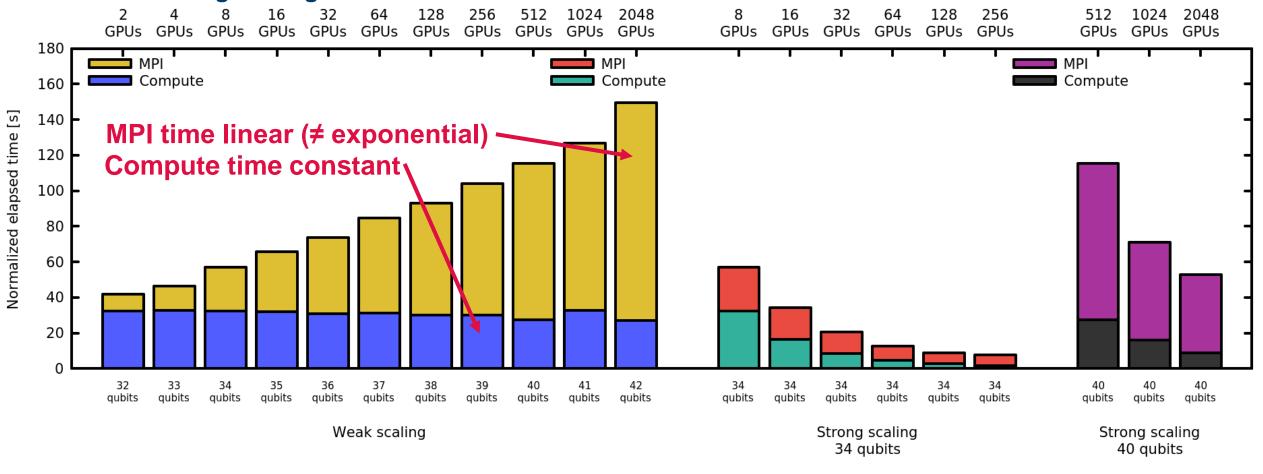




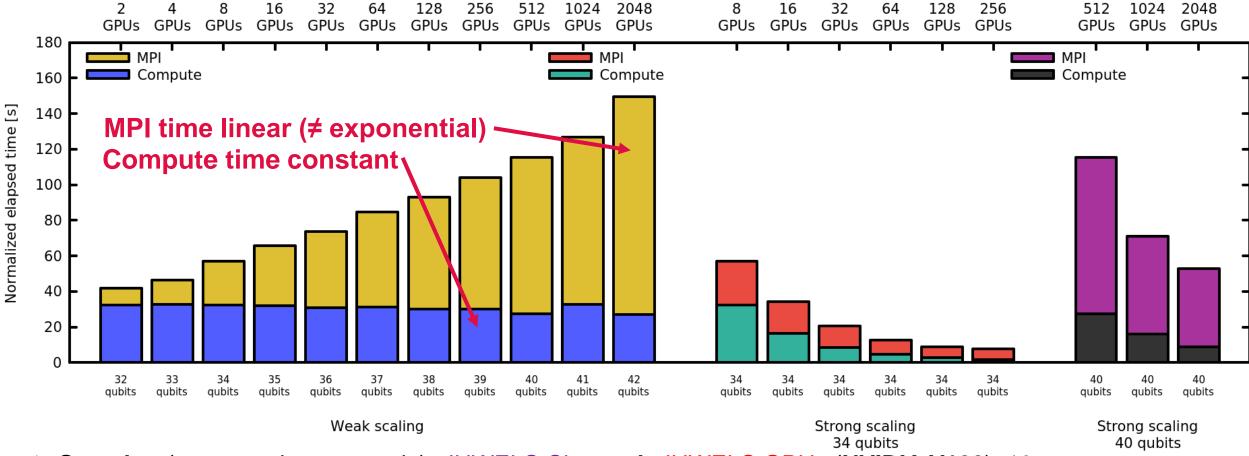






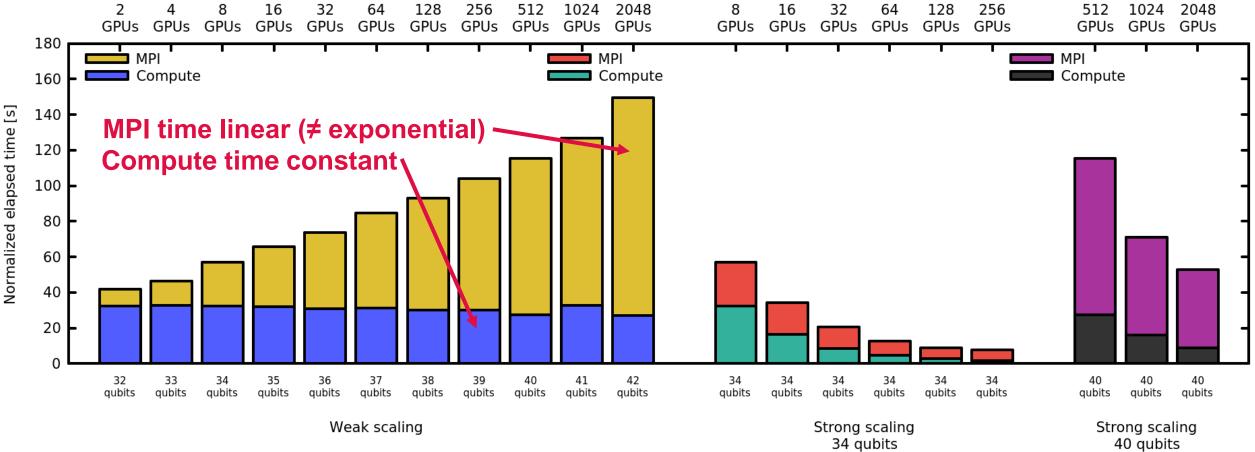


Weak and strong scaling results: Hadamard benchmark circuits



> Speedup (compute time per node): JUWELS Cluster -> JUWELS GPUs (NVIDIA V100): 10



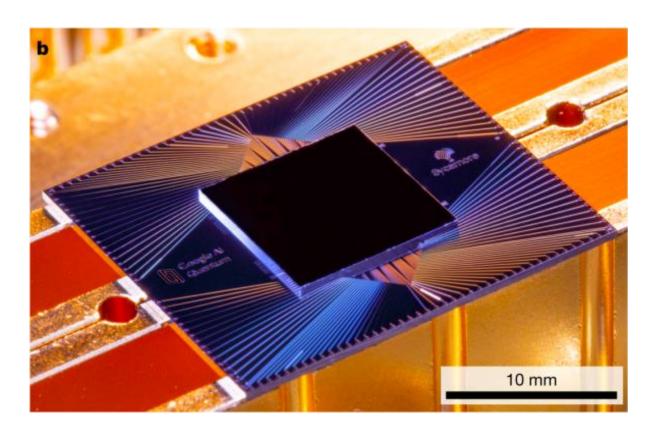


- > Speedup (compute time per node): JUWELS Cluster → JUWELS GPUs (NVIDIA V100): 10
- ➤ Speedup (compute time per node): JUWELS GPUs → JUWELS Booster (NVIDIA A100): 2 3

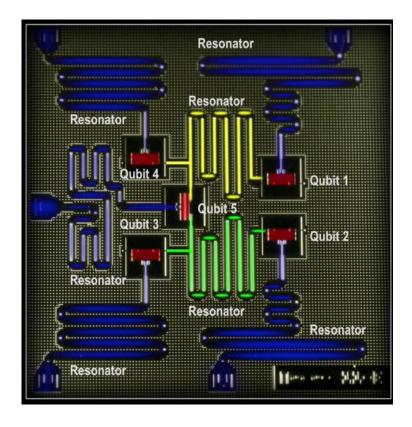


Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++







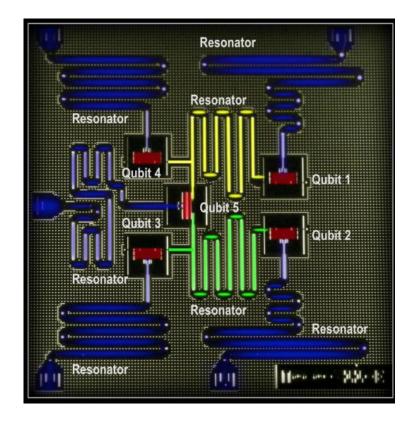


Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

> Physical realization: Solve Schrödinger / Master Equation

$$\frac{\partial}{\partial t}|\psi\rangle = -iH|\psi\rangle \quad \text{or} \quad \frac{\partial}{\partial t}\rho = -i[H,\rho] + \mathcal{D}[\rho] = \mathcal{L}[\rho]$$





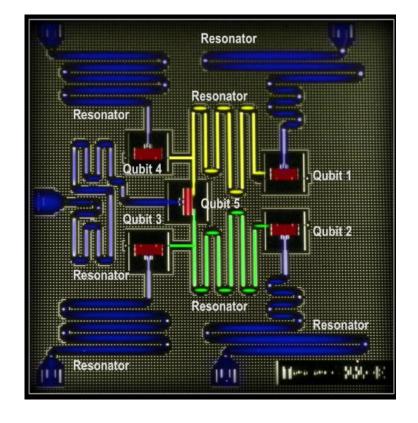
Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

➤ Physical realization: Solve Schrödinger / Master Equation

$$\frac{\partial}{\partial t}|\psi\rangle = -iH|\psi\rangle \quad \text{or} \quad \frac{\partial}{\partial t}\rho = -i[H,\rho] + \mathcal{D}[\rho] = \mathcal{L}[\rho]$$

> Similar SPMV updates **but**:





Simulating physical realizations of quantum computers on GPUs

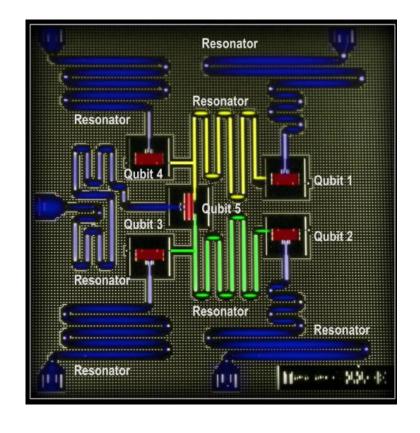
CUDA+OpenACC MPI C++

> Physical realization: Solve Schrödinger / Master Equation

$$\frac{\partial}{\partial t}|\psi\rangle = -iH|\psi\rangle \quad \text{or} \quad \frac{\partial}{\partial t}\rho = -i[H,\rho] + \mathcal{D}[\rho] = \mathcal{L}[\rho]$$

- ➤ Similar SPMV updates **but**:
 - ➤ Many updates per time step

$$\rho(t+\tau) = e^{\mathcal{L}(t+\tau/2)}\rho(t)$$





Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

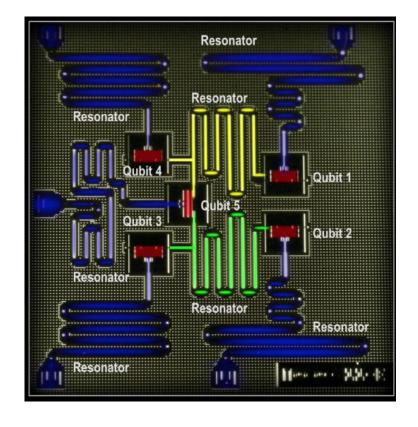
➤ Physical realization: Solve Schrödinger / Master Equation

$$\frac{\partial}{\partial t}|\psi\rangle = -iH|\psi\rangle \quad \text{or} \quad \frac{\partial}{\partial t}\rho = -i[H,\rho] + \mathcal{D}[\rho] = \mathcal{L}[\rho]$$

- ➤ Similar SPMV updates **but**:
 - ➤ Many updates per time step

$$\rho(t+\tau) = e^{\mathcal{L}(t+\tau/2)}\rho(t)$$

➤ More than 2 states per subsystem





Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

➤ Physical realization: Solve Schrödinger / Master Equation

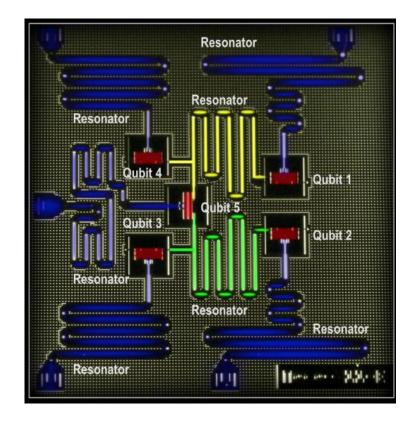
$$\frac{\partial}{\partial t}|\psi\rangle = -iH|\psi\rangle \quad \text{or} \quad \frac{\partial}{\partial t}\rho = -i[H,\rho] + \mathcal{D}[\rho] = \mathcal{L}[\rho]$$

- ➤ Similar SPMV updates **but**:
 - > Many updates per time step

$$\rho(t+\tau) = e^{\mathcal{L}(t+\tau/2)}\rho(t)$$

➤ More than 2 states per subsystem

$$ightharpoonup$$
 Before: $|\psi\rangle=(\psi_{q_3q_2q_1q_0})$





Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

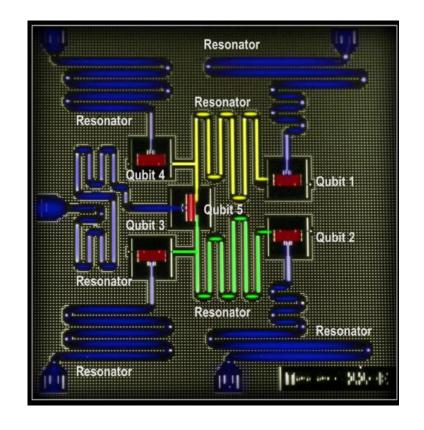
➤ Physical realization: Solve Schrödinger / Master Equation

$$\frac{\partial}{\partial t}|\psi\rangle = -iH|\psi\rangle \quad \text{or} \quad \frac{\partial}{\partial t}\rho = -i[H,\rho] + \mathcal{D}[\rho] = \mathcal{L}[\rho]$$

- ➤ Similar SPMV updates **but**:
 - > Many updates per time step

$$\rho(t+\tau) = e^{\mathcal{L}(t+\tau/2)}\rho(t)$$

- ➤ More than 2 states per subsystem
 - ightharpoonup Before: $|\psi\rangle=(\psi_{q_3q_2q_1q_0})$
 - Now: $\rho = (\rho_{k_0 m_0 k_1 m_1})$





Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

➤ Physical realization: Solve Schrödinger / Master Equation

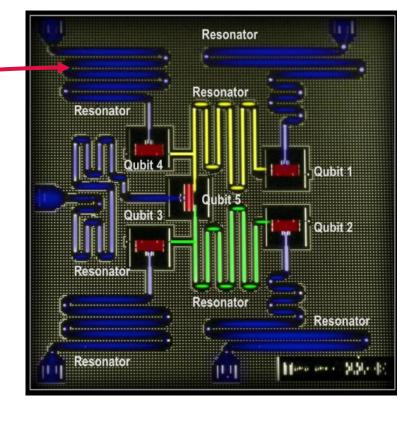
$$\frac{\partial}{\partial t}|\psi\rangle = -iH|\psi\rangle \quad \text{or} \quad \frac{\partial}{\partial t}\rho = -i[H,\rho] + \mathcal{D}[\rho] = \mathcal{L}[\rho]$$

- ➤ Similar SPMV updates **but**:
 - > Many updates per time step

$$\rho(t+\tau) = e^{\mathcal{L}(t+\tau/2)}\rho(t)$$

- ➤ More than 2 states per subsystem
 - ightharpoonup Before: $|\psi
 angle=(\psi_{q_3q_2q_1q_0})$
 - Now: $\rho = (\rho_{k_0 m_0 k_1 m_1})$

 $0 \le k_0, k_1 < 1000$





Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

> Physical realization: Solve Schrödinger / Master Equation

$$\frac{\partial}{\partial t}|\psi\rangle = -iH|\psi\rangle \quad \text{or} \quad \frac{\partial}{\partial t}\rho = -i[H,\rho] + \mathcal{D}[\rho] = \mathcal{L}[\rho]$$

- ➤ Similar SPMV updates **but**:
 - ➤ Many updates per time step

$$\rho(t+\tau) = e^{\mathcal{L}(t+\tau/2)}\rho(t)$$

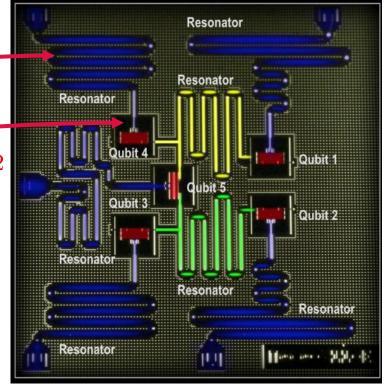
➤ More than 2 states per subsystem

$$ightharpoonup$$
 Before: $|\psi
angle=(\psi_{q_3q_2q_1q_0})$

Now:
$$\rho = (\rho_{k_0 m_0 k_1 m_1})$$

 $0 \le k_0, k_1 < 1000$

$$0 \le m_0, m_1 < 4 \text{ to } 12$$





Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

➤ Physical realization: Solve Schrödinger / Master Equation

$$\frac{\partial}{\partial t}|\psi\rangle = -iH|\psi\rangle \quad \text{or} \quad \frac{\partial}{\partial t}\rho = -i[H,\rho] + \mathcal{D}[\rho] = \mathcal{L}[\rho]$$

- ➤ Similar SPMV updates **but**:
 - > Many updates per time step

$$\rho(t+\tau) = e^{\mathcal{L}(t+\tau/2)}\rho(t)$$

➤ More than 2 states per subsystem

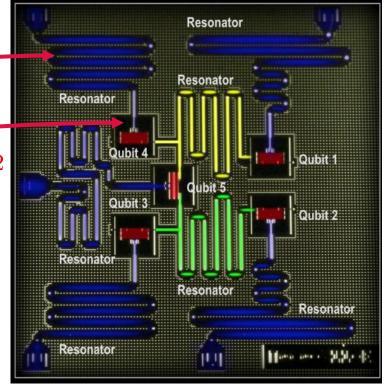
$$ightharpoonup$$
 Before: $|\psi
angle=(\psi_{q_3q_2q_1q_0})$

Now:
$$\rho = (\rho_{k_0 m_0 k_1 m_1})$$

➤ More complicated sparse matrices (sin, sinh, cos, cosh, exp, ...)

$$0 \le k_0, k_1 < 1000$$

$$0 \le m_0, m_1 < 4 \text{ to } 12$$





Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

➤ Physical realization: Solve Schrödinger / Master Equation

$$\frac{\partial}{\partial t}|\psi\rangle = -iH|\psi\rangle \quad \text{or} \quad \frac{\partial}{\partial t}\rho = -i[H,\rho] + \mathcal{D}[\rho] = \mathcal{L}[\rho]$$

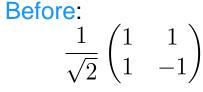
- ➤ Similar SPMV updates **but**:
 - ➤ Many updates per time step

$$\rho(t+\tau) = e^{\mathcal{L}(t+\tau/2)}\rho(t)$$

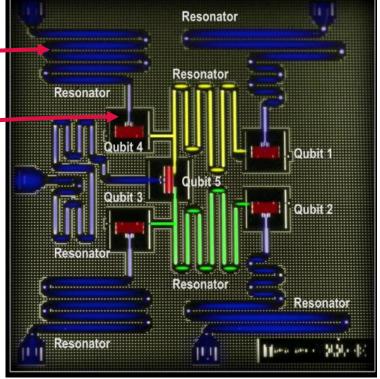
- ➤ More than 2 states per subsystem
 - ightharpoonup Before: $|\psi
 angle=(\psi_{q_3q_2q_1q_0})$
 - Now: $\rho = (\rho_{k_0 m_0 k_1 m_1})$
- ➤ More complicated sparse matrices (sin, sinh, cos, cosh, exp, ...)

$$0 \le k_0, k_1 < 1000$$

$$0 \le m_0, m_1 < 4 \text{ to } 12$$



Now:
$$\begin{pmatrix} \tilde{c}_{k_0m_0} & -i\tilde{s}_{k_0m_0} \\ -i\tilde{s}_{k_0m_0} & \tilde{c}_{k_0m_0} \end{pmatrix}$$





Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

➤ Physical realization: Solve Schrödinger / Master Equation

$$\frac{\partial}{\partial t}|\psi\rangle = -iH|\psi\rangle \quad \text{or} \quad \frac{\partial}{\partial t}\rho = -i[H,\rho] + \mathcal{D}[\rho] = \mathcal{L}[\rho]$$

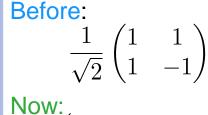
- ➤ Similar SPMV updates **but**:
 - ➤ Many updates per time step

$$\rho(t+\tau) = e^{\mathcal{L}(t+\tau/2)}\rho(t)$$

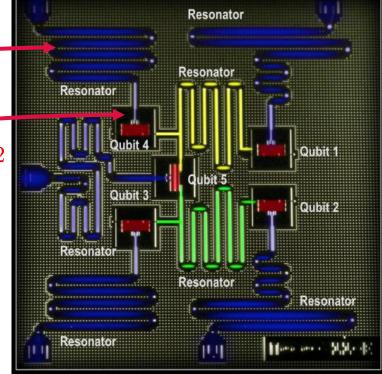
- ➤ More than 2 states per subsystem
 - ightharpoonup Before: $|\psi
 angle=(\psi_{q_3q_2q_1q_0})$
 - Now: $\rho = (\rho_{k_0 m_0 k_1 m_1})$
- ➤ More complicated sparse matrices (sin, sinh, cos, cosh, exp, ...)

$$0 \le k_0, k_1 < 1000$$

$$0 \le m_0, m_1 < 4 \text{ to } 12$$



$$egin{pmatrix} ilde{c}_{k_0m_0} & -i ilde{s}_{k_0m_0} \ -i ilde{s}_{k_0m_0} & ilde{c}_{k_0m_0} \end{pmatrix}$$



$$\tilde{c}_{k_0 m_0} = \cos(\tau \sqrt{k_0 + 1} (G\lambda_{m_0} + \varepsilon(\tilde{t})))$$

$$\tilde{s}_{k_0 m_0} = \sin(\tau \sqrt{k_0 + 1} (G\lambda_{m_0} + \varepsilon(\tilde{t})))$$



Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

➤ Physical realization: Solve Schrödinger / Master Equation

$$\frac{\partial}{\partial t}|\psi\rangle = -iH|\psi\rangle \quad \text{or} \quad \frac{\partial}{\partial t}\rho = -i[H,\rho] + \mathcal{D}[\rho] = \mathcal{L}[\rho]$$

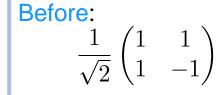
- ➤ Similar SPMV updates **but**:
 - Many updates per time step

$$\rho(t+\tau) = e^{\mathcal{L}(t+\tau/2)}\rho(t)$$

- ➤ More than 2 states per subsystem
 - ightharpoonup Before: $|\psi
 angle=(\psi_{q_3q_2q_1q_0})$
 - Now: $\rho = (\rho_{k_0 m_0 k_1 m_1})$
- ➤ More complicated sparse matrices (sin, sinh, cos, cosh, exp, ...)
- Very computation-intensive (memory "only" 2 GiB)

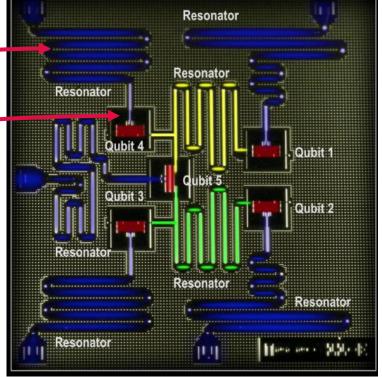
$$0 \le k_0, k_1 < 1000$$

$$0 \le m_0, m_1 < 4 \text{ to } 12$$



Page 22

Now:



$$\tilde{c}_{k_0 m_0} = \cos(\tau \sqrt{k_0 + 1} (G\lambda_{m_0} + \varepsilon(\tilde{t})))$$

$$\tilde{s}_{k_0 m_0} = \sin(\tau \sqrt{k_0 + 1} (G\lambda_{m_0} + \varepsilon(\tilde{t})))$$

$$\tilde{s}_{k_0 m_0} = \sin(\tau \sqrt{k_0 + 1} (G\lambda_{m_0} + \varepsilon(\tilde{t})))$$



Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

➤ Physical realization: Solve Schrödinger / Master Equation

$$\frac{\partial}{\partial t}|\psi\rangle = -iH|\psi\rangle \quad \text{or} \quad \frac{\partial}{\partial t}\rho = -i[H,\rho] + \mathcal{D}[\rho] = \mathcal{L}[\rho]$$

- ➤ Similar SPMV updates **but**:
 - Many updates per time step

$$\rho(t+\tau) = e^{\mathcal{L}(t+\tau/2)}\rho(t)$$

- ➤ More than 2 states per subsystem
 - ightharpoonup Before: $|\psi
 angle=(\psi_{q_3q_2q_1q_0})$
 - Now: $\rho = (\rho_{k_0 m_0 k_1 m_1})$
- ➤ More complicated sparse matrices (sin, sinh, cos, cosh, exp, ...)
- Very computation-intensive (memory "only" 2 GiB)
 - → Useful to measure single-GPU performance

$$0 \le k_0, k_1 < 1000$$

$$0 \le m_0, m_1 < 4 \text{ to } 12$$

Before: $\frac{1}{\sqrt{2}}\begin{pmatrix}1&1\\1&-1\end{pmatrix}$

Page 22

Now:

$$\tilde{c}_{k_0 m_0} = \cos(\tau \sqrt{k_0 + 1} (G\lambda_{m_0} + \varepsilon(\tilde{t})))$$

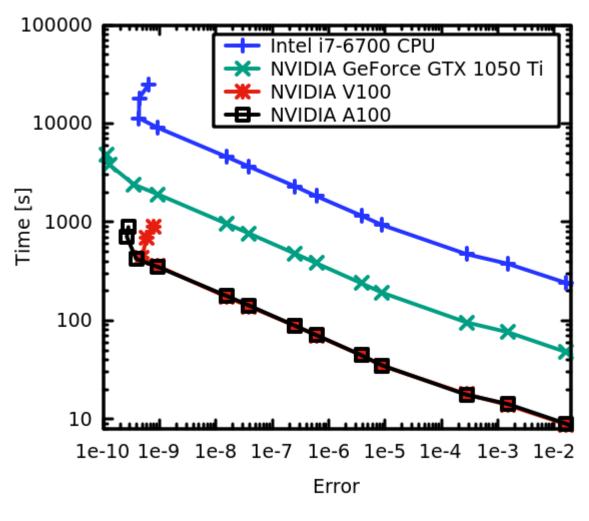
$$\tilde{s}_{k_0 m_0} = \sin(\tau \sqrt{k_0 + 1} (G\lambda_{m_0} + \varepsilon(\tilde{t})))$$

$$\tilde{s}_{k_0 m_0} = \sin(\tau \sqrt{k_0 + 1} (G\lambda_{m_0} + \varepsilon(\tilde{t})))$$



Simulating physical realizations of quantum computers on GPUs

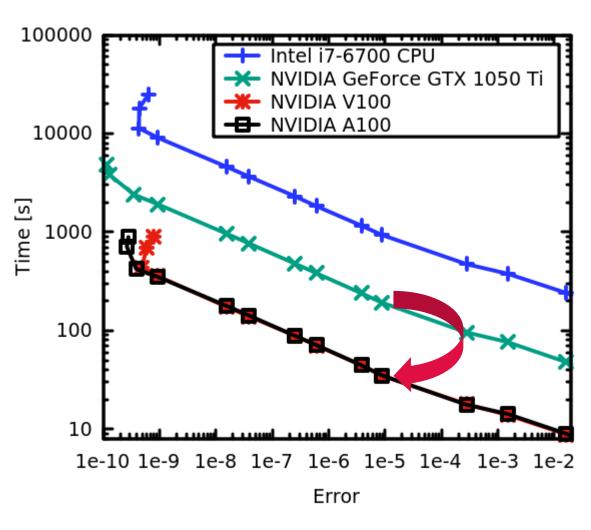
CUDA+OpenACC MPI C++





Simulating physical realizations of quantum computers on GPUs

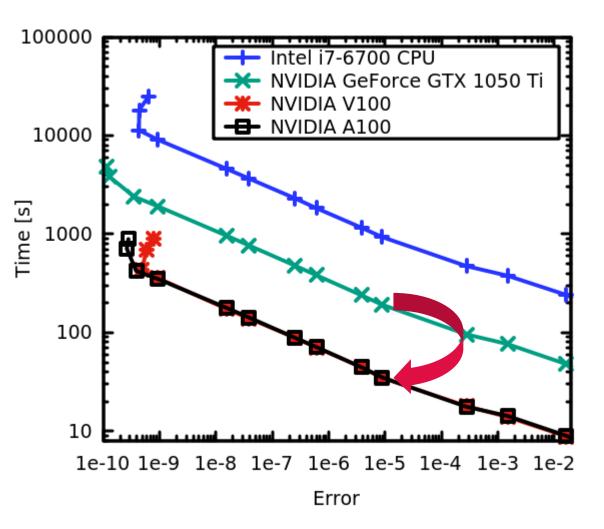
CUDA+OpenACC MPI C++





Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++



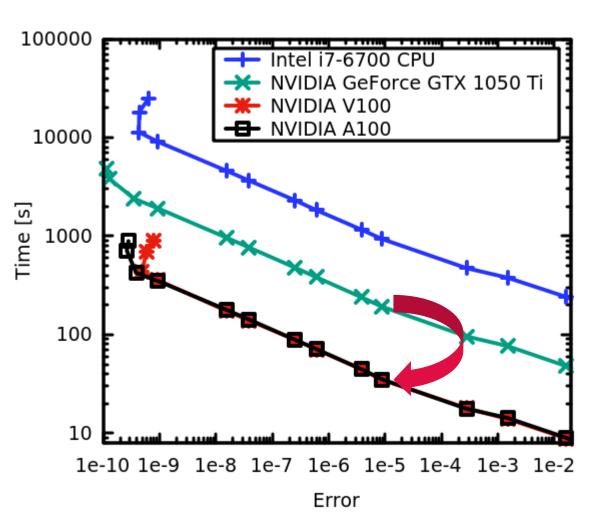
For each
$$k_0$$
 in steps of 2, for each k_1, m_0, m_1

$$\begin{pmatrix} \rho_{k_0 m_0 k_1 m_1} \\ \rho_{(k_0+1)m_0(k_1+1)m_1} \end{pmatrix} \leftarrow \begin{pmatrix} c_{k_0, k_1} & s_{k_0, k_1}^- \\ s_{k_0, k_1}^+ & c_{k_0, k_1} \end{pmatrix} \begin{pmatrix} \rho_{k_0 m_0 k_1 m_1} \\ \rho_{(k_0+1)m_0(k_1+1)m_1} \end{pmatrix}$$



Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

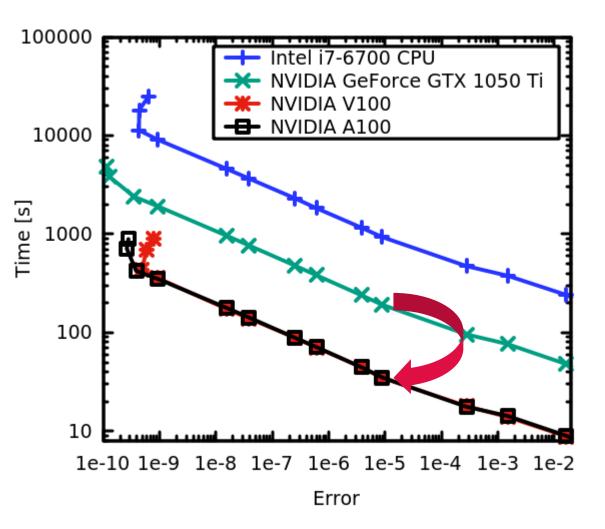


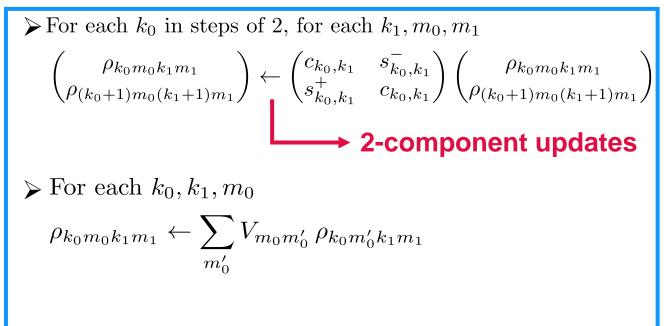


Page 23

Simulating physical realizations of quantum computers on GPUs

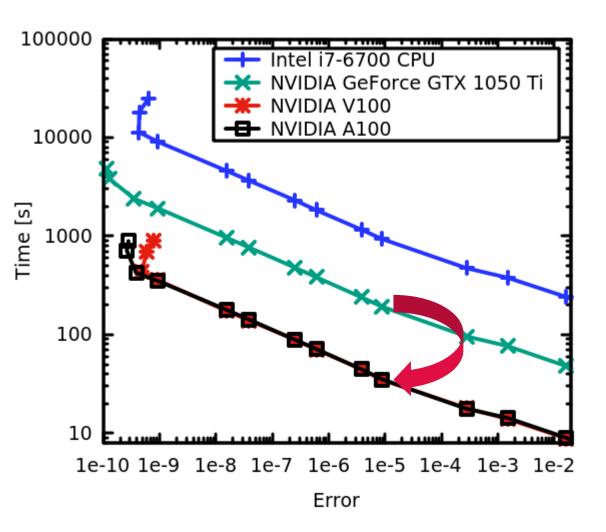
CUDA+OpenACC MPI C++

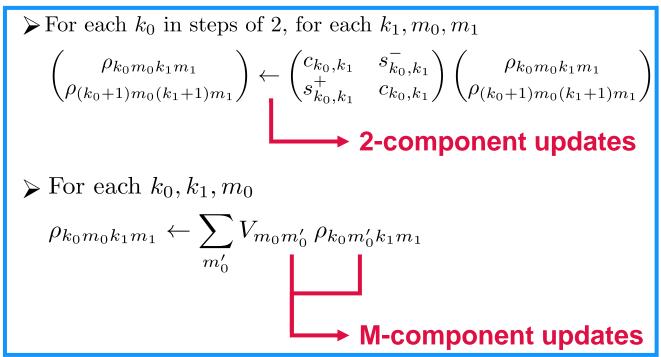




Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

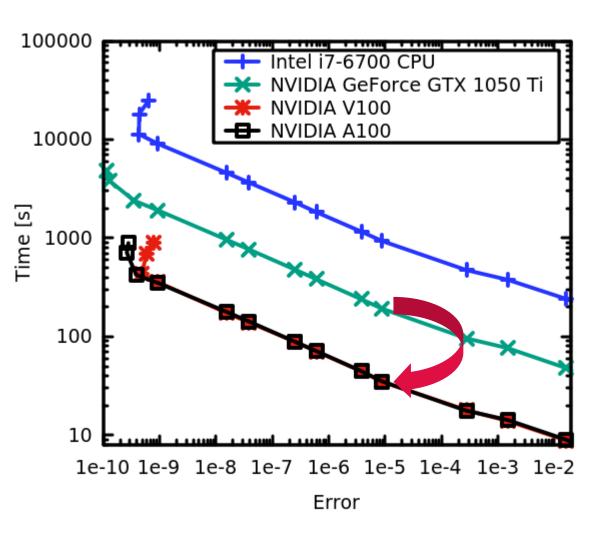




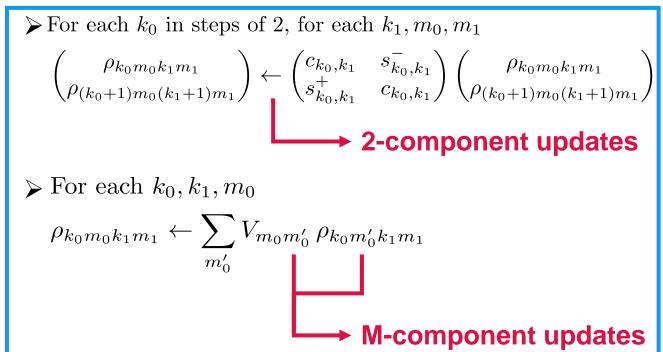


Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++



Bottleneck kernels profit from **Tensor Cores**

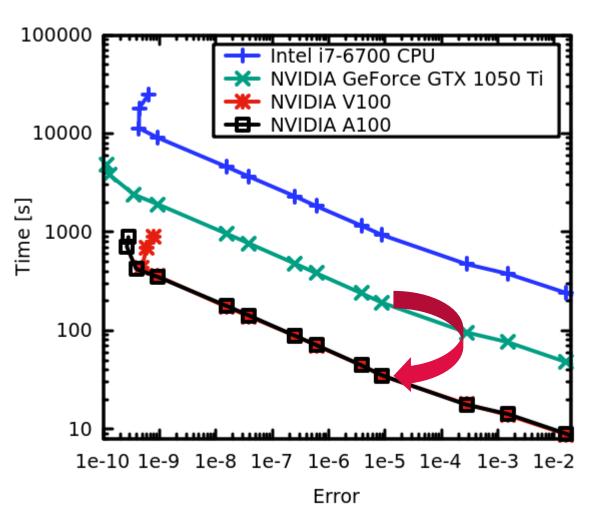


Small system with M=4 and K=100

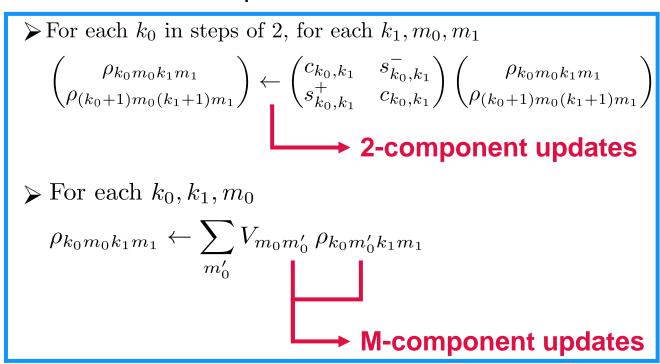


Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++



Bottleneck kernels profit from **Tensor Cores**



Small system with M=4 and K=100

→ Similar performance for V100 and A100



Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

Realistic system with M=12 and K=400



Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

Realistic system with M=12 and K=400

For each
$$k_0, k_1, m_0$$

$$\rho_{k_0 m_0 k_1 m_1} \leftarrow \sum_{m_0'} V_{m_0 m_0'} \, \rho_{k_0 m_0' k_1 m_1}$$

12-component updates



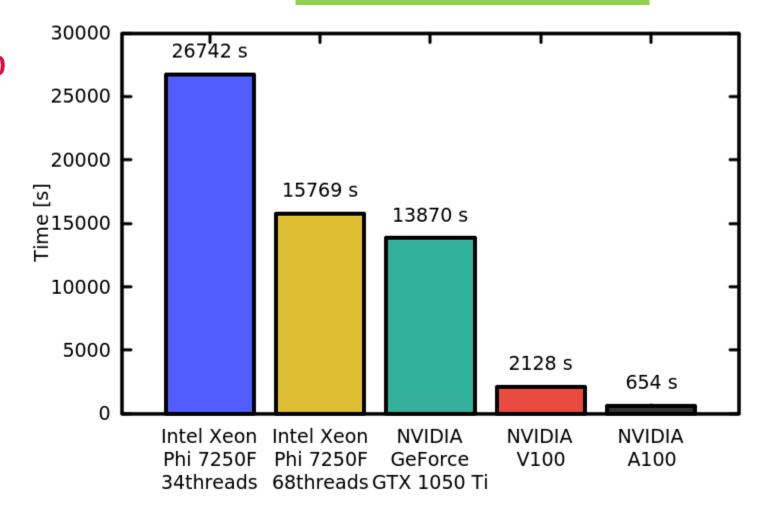
Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

Realistic system with M=12 and K=400

For each k_0, k_1, m_0 $\rho_{k_0 m_0 k_1 m_1} \leftarrow \sum_{m'_0} V_{m_0 m'_0} \rho_{k_0 m'_0 k_1 m_1}$

12-component updates





Page 24

Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

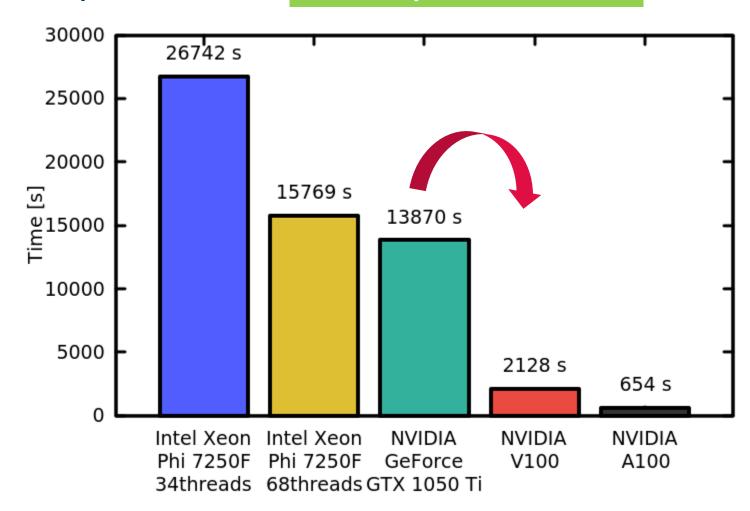
Realistic system with M=12 and K=400

For each
$$k_0, k_1, m_0$$

$$\rho_{k_0 m_0 k_1 m_1} \leftarrow \sum_{m'_0} V_{m_0 m'_0} \rho_{k_0 m'_0 k_1 m_1}$$

12-component updates

→ Similar speedup until **V100**





Page 24

Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++

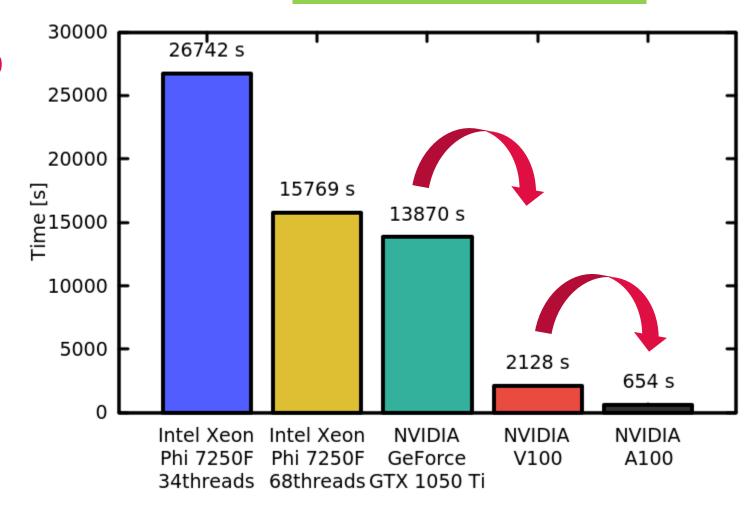
Realistic system with M=12 and K=400

For each
$$k_0, k_1, m_0$$

$$\rho_{k_0 m_0 k_1 m_1} \leftarrow \sum_{m'_0} V_{m_0 m'_0} \rho_{k_0 m'_0 k_1 m_1}$$

12-component updates

- → Similar speedup until **V100**
- → Further speedup from **V100** to **A100**









> All QC simulations are composed of



sparse

complex

matrix-vector updates





JÜLICH SUPERCOMPUTING CENTRE

> All QC simulations are composed of





complex

matrix-vector updates

> Simulating QCs is a versatile approach to benchmark supercomputers



JÜLICH SUPERCOMPUTING CENTRE

> All QC simulations are composed of





complex

matrix-vector updates

- > Simulating QCs is a versatile approach to benchmark supercomputers
 - ➤ Memory-, network-, and computation-intensive



JÜLICH SUPERCOMPUTING CENTRE

> All QC simulations are composed of







matrix-vector updates

- > Simulating QCs is a versatile approach to benchmark supercomputers
 - > Memory-, network-, and computation-intensive
- > Huge speedup from CPU-based simulators to GPU-based simulators



JÜLICH SUPERCOMPUTING CENTRE

> All QC simulations are composed of







matrix-vector updates

- > Simulating QCs is a versatile approach to benchmark supercomputers
 - > Memory-, network-, and computation-intensive
- > Huge speedup from CPU-based simulators to GPU-based simulators

THANK YOU FOR YOUR ATTENTION

- More information and references:
 - > MPI communication scheme: De Raedt et al., Comp. Phys. Commun. 176, 121 (2007)
 - > Benchmarking gate-based quantum computers: Michielsen et al., Comp. Phys. Commun. 220, 44 (2017)
 - > JUQCS: De Raedt et al., Comp. Phys. Commun. 237, 41 (2019)
 - Quantum supremacy with JUQCS: Arute et al., Nature 574, 505 (2019)
 - Benchmarking supercomputers with JUQCS: Willsch et al., NIC Series 50, 255 (2020)
 - GPU-accelerated simulations of QA and the QAOA: Willsch et al., preprint available on arXiv (2021)



What are the fundamental tensor operations in all simulations of this kind?



What are the fundamental tensor operations in all simulations of this kind?

$$\psi_{\cdots q_{3}q_{2}q_{1}q_{0}} \leftarrow \sum_{q'_{2}} H_{q_{2}q'_{2}} \psi_{\cdots q_{3}q'_{2}q_{1}q_{0}}$$

$$\psi_{\cdots q_{3}q_{2}q_{1}q_{0}} \leftarrow \sum_{q'_{2}q'_{0}} U_{q_{2}q'_{2}q_{0}q'_{0}} \psi_{\cdots q_{3}q'_{2}q_{1}q'_{0}}$$

$$\rho_{k_{0}m_{0}k_{1}m_{1}\cdots} \leftarrow \sum_{m'_{0}} V_{m_{0}m'_{0}} \rho_{k_{0}m'_{0}k_{1}m_{1}\cdots}$$

$$\rho_{k_{0}m_{0}k_{1}m_{1}\cdots} \leftarrow \sum_{k'_{0}k'_{1}} W_{k_{0}k'_{0}k_{1}k'_{1}} \rho_{k'_{0}m_{0}k'_{1}m_{1}\cdots}$$



What are the fundamental tensor operations in all simulations of this kind?

$$\psi_{\dots q_3 q_2 q_1 q_0} \leftarrow \sum_{q_2'} H_{q_2 q_2'} \psi_{\dots q_3 q_2' q_1 q_0}$$

$$\psi_{\dots q_3 q_2 q_1 q_0} \leftarrow \sum_{q_2' q_0'} U_{q_2 q_2' q_0 q_0'} \psi_{\dots q_3 q_2' q_1 q_0'}$$

$$\rho_{k_0 m_0 k_1 m_1 \dots} \leftarrow \sum_{m_0'} V_{m_0 m_0'} \rho_{k_0 m_0' k_1 m_1 \dots}$$

$$\rho_{k_0 m_0 k_1 m_1 \dots} \leftarrow \sum_{k_0' k_1'} W_{k_0 k_0' k_1 k_1'} \rho_{k_0' m_0 k_1' m_1 \dots}$$

What are the fundamental tensor operations in all simulations of this kind?

$$\psi_{\cdots q_{3}q_{2}q_{1}q_{0}} \leftarrow \sum_{q'_{2}} H_{q_{2}q'_{2}} \psi_{\cdots q_{3}q'_{2}q_{1}q_{0}}$$

$$\psi_{\cdots q_{3}q_{2}q_{1}q_{0}} \leftarrow \sum_{q'_{2}q'_{0}} U_{q_{2}q'_{2}q_{0}q'_{0}} \psi_{\cdots q_{3}q'_{2}q_{1}q'_{0}}$$

$$\rho_{k_{0}m_{0}k_{1}m_{1}\cdots} \leftarrow \sum_{m'_{0}} V_{m_{0}m'_{0}} \rho_{k_{0}m'_{0}k_{1}m_{1}\cdots}$$

$$\rho_{k_{0}m_{0}k_{1}m_{1}\cdots} \leftarrow \sum_{k'_{0}k'_{1}} W_{k_{0}k'_{0}k_{1}k'_{1}} \rho_{k'_{0}m_{0}k'_{1}m_{1}\cdots}$$

$$\psi_{\cdots q_{3}q_{2}q_{1}q_{0}} \leftarrow \sum_{q'_{2}} H_{q_{2}q'_{2}} \psi_{\cdots q_{3}q'_{2}q_{1}q_{0}}$$

$$\psi_{\cdots q_{3}q_{2}q_{1}q_{0}} \leftarrow \sum_{q'_{2}q'_{0}} U_{q_{2}q'_{2}q_{0}q'_{0}} \psi_{\cdots q_{3}q'_{2}q_{1}q'_{0}}$$

$$\rho_{k_{0}m_{0}k_{1}m_{1}\cdots} \leftarrow \sum_{m'_{0}} V_{m_{0}m'_{0}} \rho_{k_{0}m'_{0}k_{1}m_{1}\cdots}$$

$$For each k_{0} in steps of 2, for each k_{1}, m_{0}, m_{1}$$

$$For each k_{0} in steps of 2, for each k_{1}, m_{0}, m_{1}$$

$$\begin{pmatrix} \rho_{k_0 m_0 k_1 m_1} \\ \rho_{(k_0+1) m_0 (k_1+1) m_1} \end{pmatrix} \leftarrow \begin{pmatrix} c_{k_0, k_1} & s_{k_0, k_1}^- \\ s_{k_0, k_1}^+ & c_{k_0, k_1} \end{pmatrix} \begin{pmatrix} \rho_{k_0 m_0 k_1 m_1} \\ \rho_{(k_0+1) m_0 (k_1+1) m_1} \end{pmatrix}$$

What are the fundamental tensor operations in all simulations of this kind?

Page 26

$$\psi_{\dots q_3 q_2 q_1 q_0} \leftarrow \sum_{q_2'} H_{q_2 q_2'} \psi_{\dots q_3 q_2' q_1 q_0}$$

$$\psi_{\dots q_3 q_2 q_1 q_0} \leftarrow \sum_{q_2' q_0'} U_{q_2 q_2' q_0 q_0'} \psi_{\dots q_3 q_2' q_1 q_0'}$$

$$\rho_{k_0 m_0 k_1 m_1 \dots} \leftarrow \sum_{m_0'} V_{m_0 m_0'} \rho_{k_0 m_0' k_1 m_1 \dots}$$

$$\rho_{k_0 m_0 k_1 m_1 \dots} \leftarrow \sum_{k_0' k_1'} W_{k_0 k_0' k_1 k_1'} \rho_{k_0' m_0 k_1' m_1 \dots}$$

$$\begin{array}{c} \psi_{\cdots q_{3}q_{2}q_{1}q_{0}} \leftarrow \sum_{q_{2}'} H_{q_{2}q_{2}'} \psi_{\cdots q_{3}q_{2}'q_{1}q_{0}} \\ \psi_{\cdots q_{3}q_{2}q_{1}q_{0}} \leftarrow \sum_{q_{2}'q_{0}'} U_{q_{2}q_{2}'q_{0}q_{0}'} \psi_{\cdots q_{3}q_{2}'q_{1}q_{0}'} \\ \rho_{k_{0}m_{0}k_{1}m_{1}\cdots} \leftarrow \sum_{m_{0}'} V_{m_{0}m_{0}'} \rho_{k_{0}m_{0}'k_{1}m_{1}\cdots} \end{array} \qquad \begin{array}{c} \sum_{k_{1}'} \begin{pmatrix} 1 \\ c_{11} & -is_{12} \\ -is_{21} & c_{22} \\ & \ddots \\ & c_{(K-2)(K-2)} & -is_{(K-2)(K-1)} \\ & -is_{(K-2)(K-1)} & c_{(K-1)(K-1)} \end{pmatrix}_{k_{1}k_{1}'} \\ \rho_{k_{0}m_{0}k_{1}'m_{1}\cdots} \\ & \sum_{m_{0}'} W_{n_{0}m_{0}'} \rho_{k_{0}m_{0}'k_{1}m_{1}\cdots} \end{array}$$

For each k_0 in steps of 2, for each k_1, m_0, m_1

$$\begin{pmatrix} \rho_{k_0 m_0 k_1 m_1} \\ \rho_{(k_0+1) m_0 (k_1+1) m_1} \end{pmatrix} \leftarrow \begin{pmatrix} c_{k_0, k_1} & s_{k_0, k_1}^- \\ s_{k_0, k_1}^+ & c_{k_0, k_1} \end{pmatrix} \begin{pmatrix} \rho_{k_0 m_0 k_1 m_1} \\ \rho_{(k_0+1) m_0 (k_1+1) m_1} \end{pmatrix}$$

→ In-place double complex 2D SPMV updates



What dedicated library functions would be useful for quantum computer simulations?



What dedicated library functions would be useful for quantum computer simulations?

$$\psi..._{q_3q_2q_1q_0} \leftarrow \sum_{q_2'} H_{q_2q_2'} \psi..._{q_3q_2'q_1q_0}$$

$$\psi..._{q_3q_2q_1q_0} \leftarrow \sum_{q_2'q_0'} U_{q_2q_2'q_0q_0'} \psi..._{q_3q_2'q_1q_0'}$$

$$\rho_{k_0m_0k_1m_1}... \leftarrow \sum_{m_0'} V_{m_0m_0'} \rho_{k_0m_0'k_1m_1}...$$

$$\rho_{k_0m_0k_1m_1}... \leftarrow \sum_{k_0'k_1'} W_{k_0k_0'k_1k_1'} \rho_{k_0'm_0k_1'm_1}...$$

What dedicated library functions would be useful for quantum computer simulations?

$$\psi_{\cdots q_{3}q_{2}q_{1}q_{0}} \leftarrow \sum_{q'_{2}} H_{q_{2}q'_{2}} \psi_{\cdots q_{3}q'_{2}q_{1}q_{0}}$$

$$\psi_{\cdots q_{3}q_{2}q_{1}q_{0}} \leftarrow \sum_{q'_{2}q'_{0}} U_{q_{2}q'_{2}q_{0}q'_{0}} \psi_{\cdots q_{3}q'_{2}q_{1}q'_{0}}$$

$$\rho_{k_{0}m_{0}k_{1}m_{1}\cdots} \leftarrow \sum_{m'_{0}} V_{m_{0}m'_{0}} \rho_{k_{0}m'_{0}k_{1}m_{1}\cdots}$$

$$\rho_{k_{0}m_{0}k_{1}m_{1}\cdots} \leftarrow \sum_{k'_{0}k'_{1}} W_{k_{0}k'_{0}k_{1}k'_{1}} \rho_{k'_{0}m_{0}k'_{1}m_{1}\cdots}$$

Generic 1-component and 2-component in-place double complex 2D SPMV updates



What dedicated library functions would be useful for quantum computer simulations?

$$\psi_{\cdots q_{3}q_{2}q_{1}q_{0}} \leftarrow \sum_{q'_{2}} H_{q_{2}q'_{2}} \psi_{\cdots q_{3}q'_{2}q_{1}q_{0}}$$

$$\psi_{\cdots q_{3}q_{2}q_{1}q_{0}} \leftarrow \sum_{q'_{2}q'_{0}} U_{q_{2}q'_{2}q_{0}q'_{0}} \psi_{\cdots q_{3}q'_{2}q_{1}q'_{0}}$$

$$\rho_{k_{0}m_{0}k_{1}m_{1}\cdots} \leftarrow \sum_{m'_{0}} V_{m_{0}m'_{0}} \rho_{k_{0}m'_{0}k_{1}m_{1}\cdots}$$

$$\rho_{k_{0}m_{0}k_{1}m_{1}\cdots} \leftarrow \sum_{k'_{0}k'_{1}} W_{k_{0}k'_{0}k_{1}k'_{1}} \rho_{k'_{0}m_{0}k'_{1}m_{1}\cdots}$$

Generic 1-component and 2-component in-place double complex 2D SPMV updates

What dedicated library functions would be useful for quantum computer simulations?

$$\psi..._{q_{3}q_{2}q_{1}q_{0}} \leftarrow \sum_{q'_{2}} H_{q_{2}q'_{2}} \psi..._{q_{3}q'_{2}q_{1}q_{0}}$$

$$\psi..._{q_{3}q_{2}q_{1}q_{0}} \leftarrow \sum_{q'_{2}q'_{0}} U_{q_{2}q'_{2}q_{0}q'_{0}} \psi..._{q_{3}q'_{2}q_{1}q'_{0}}$$

$$\rho_{k_{0}m_{0}k_{1}m_{1}} \cdots \leftarrow \sum_{m'_{0}} V_{m_{0}m'_{0}} \rho_{k_{0}m'_{0}k_{1}m_{1}} \cdots$$

$$\rho_{k_{0}m_{0}k_{1}m_{1}} \cdots \leftarrow \sum_{k'_{0}k'_{1}} W_{k_{0}k'_{0}k_{1}k'_{1}} \rho_{k'_{0}m_{0}k'_{1}m_{1}} \cdots$$

Generic 1-component and 2-component in-place double complex 2D SPMV updates



Do you make use of cuBLAS or other dedicated libraries in your code?



Do you make use of cuBLAS or other dedicated libraries in your code?

Page 28

Yes, but:



Do you make use of cuBLAS or other dedicated libraries in your code?

Yes, but:

> cuBLAS:



Do you make use of cuBLAS or other dedicated libraries in your code?

Yes, but:

> cuBLAS:

$$C = \alpha \operatorname{op}(A)\operatorname{op}(B) + \beta C$$



Do you make use of cuBLAS or other dedicated libraries in your code?

Yes, but:

> cuBLAS:

cublasZgemm:

$$C = \alpha \operatorname{op}(A)\operatorname{op}(B) + \beta C$$

cublasGemmStridedBatchedEx:

$$C + i * strideC = \alpha op(A + i * strideA)op(B + i * strideB) + \beta(C + i * strideC)$$

Do you make use of cuBLAS or other dedicated libraries in your code?

Yes, but:

> cuBLAS:

cublasZgemm:

$$C = \alpha \operatorname{op}(A)\operatorname{op}(B) + \beta C$$

cublasGemmStridedBatchedEx:

$$C + i * strideC = \alpha op(A + i * strideA)op(B + i * strideB) + \beta(C + i * strideC)$$

> cuTENSOR: possible, but typically support for in-place operations missing



Do you make use of cuBLAS or other dedicated libraries in your code?

Yes, but:

> cuBLAS:

cublasZgemm: $C = \alpha op(A)op(B) + \beta C$

cublasGemmStridedBatchedEx:

$$C + i * \text{strideC} = \alpha \text{op}(A + i * \text{strideA}) \text{op}(B + i * \text{strideB}) + \beta(C + i * \text{strideC})$$

- > cuTENSOR: possible, but typically support for in-place operations missing
- > Other libraries: often support for <double> times <double complex> missing



Do you make use of cuBLAS or other dedicated libraries in your code?

Yes, but:

> cuBLAS:

cublasZgemm:

$$C = \alpha \operatorname{op}(A)\operatorname{op}(B) + \beta C$$

Page 28

cublasGemmStridedBatchedEx:

$$C + i * \text{strideC} = \alpha \text{op}(A + i * \text{strideA}) \text{op}(B + i * \text{strideB}) + \beta(C + i * \text{strideC})$$

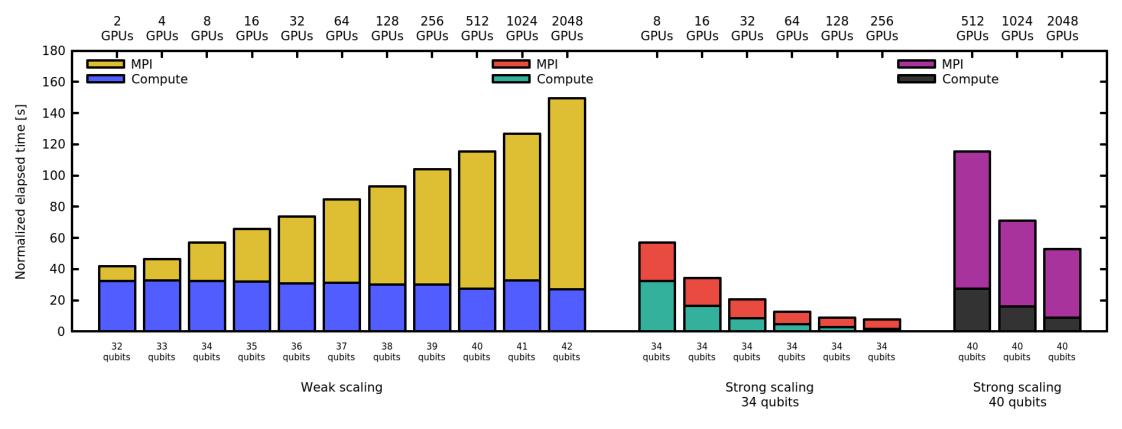
- > cutensor: possible, but typically support for in-place operations missing
- Other libraries: often support for <double> times <double complex> missing
- Our custom CUDA kernels were faster



What is the bottleneck that limits the performance of JUQCS and are there ideas to improve it?



What is the bottleneck that limits the performance of JUQCS and are there ideas to improve it?

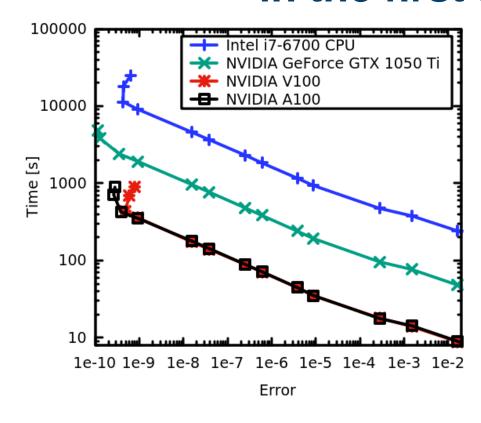




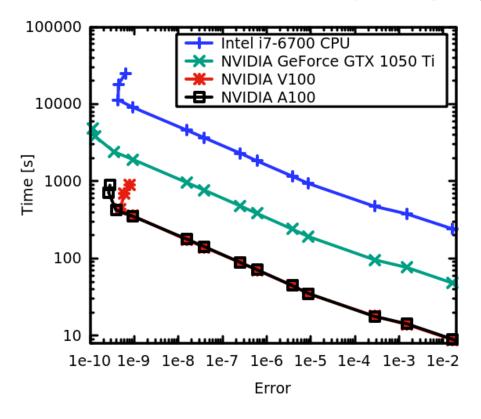
Why did V100 and A100 perform equally well in the first JUQMES result?

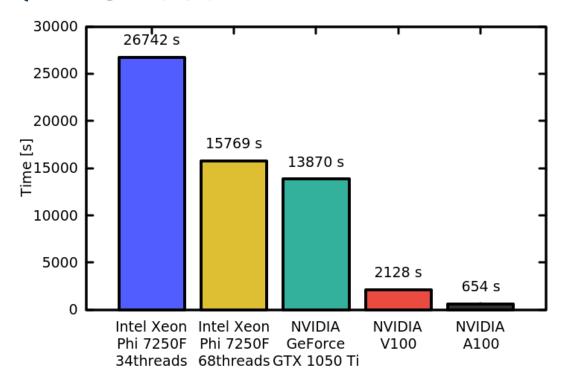


Why did V100 and A100 perform equally well in the first JUQMES result?



Why did V100 and A100 perform equally well in the first JUQMES result?

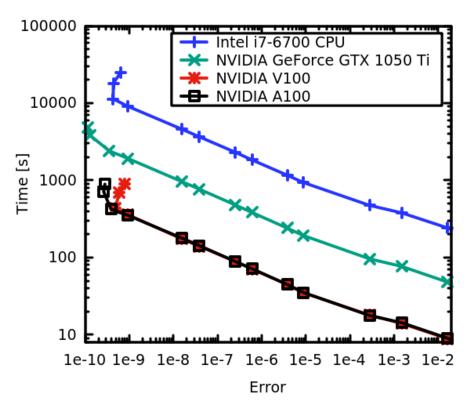


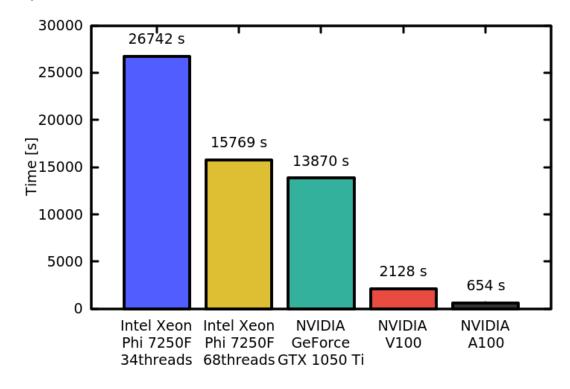




$$\rho = (\rho_{k_0 m_0 k_1 m_1}) \longrightarrow \text{Memory: } K \times M \times K \times M$$

Why did V100 and A100 perform equally well in the first JUQMES result?

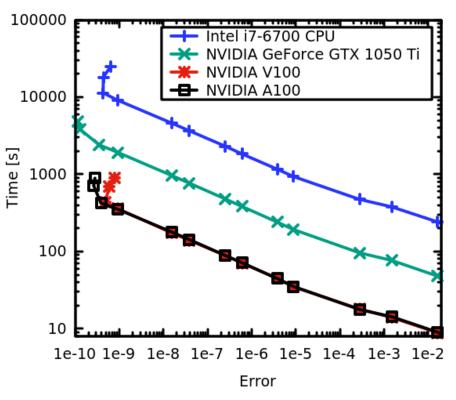


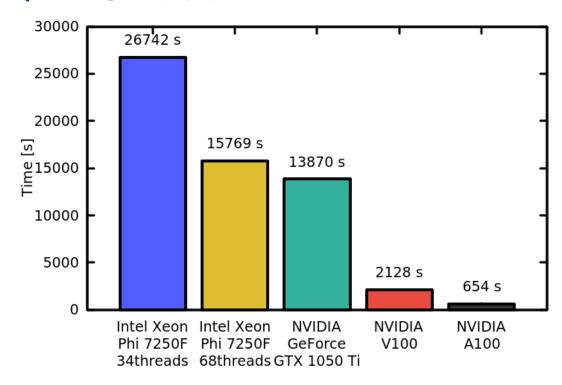




$$\rho = (\rho_{k_0 m_0 k_1 m_1}) \longrightarrow \text{Memory: } K \times M \times K \times M$$

Why did V100 and A100 perform equally well in the first JUQMES result?





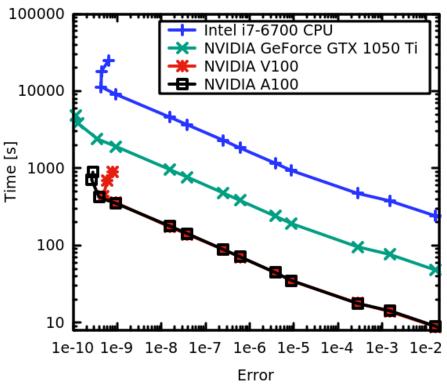
Small system with M=4 and K=100



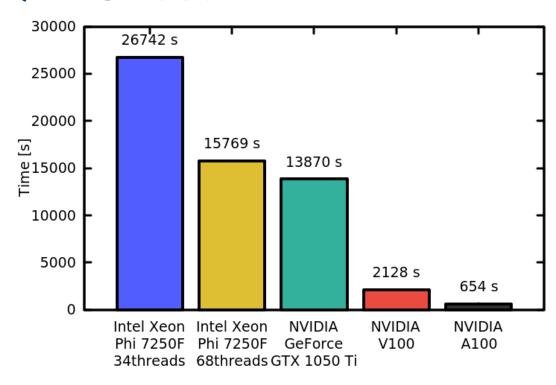
$$\rho = (\rho_{k_0 m_0 k_1 m_1}) \longrightarrow \text{Memory: } K \times M \times K \times M$$

Why did V100 and A100 perform equally well in the first JUQMES result?

Page 30



Small system with M=4 and K=100



Realistic system with M=12 and K=400

