

BENCHMARKING GPU CLUSTERS WITH THE JÜLICH UNIVERSAL QUANTUM COMPUTER SIMULATOR

JUNE 24 I 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

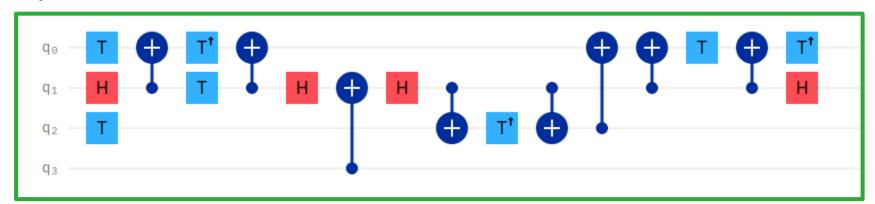


Prof. Dr. Kristel **Michielsen**



Ideal gate-based quantum computing

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



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



complex



Willsch et al. (2021

arXiv:2104.03293

> with huge vectors and huge² matrices



Ideal gate-based quantum computing

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

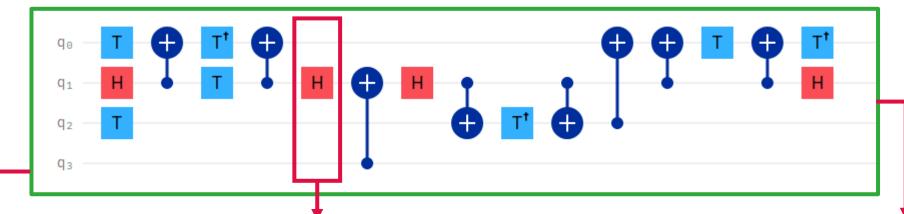
$$|\psi\rangle = \psi_{0\cdots 0}|0\cdots 0\rangle + \cdots + \psi_{1\cdots 1}|1\cdots 1\rangle = \begin{pmatrix} \psi_{0\cdots 0} \\ \vdots \\ \psi_{1\cdots 1} \end{pmatrix}$$

> 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:

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

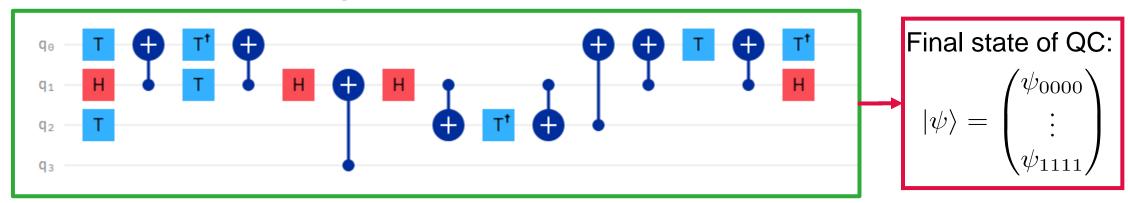
Example matrix-vector multiplication: \blacksquare on qubit q_1

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

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?
 - > 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 $q_{n-1}\cdots q_1q_0$

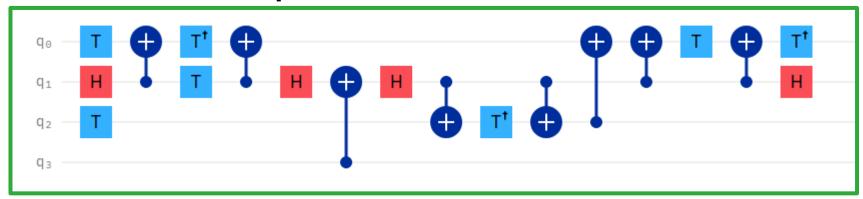
Page 5

> Need to run circuit repeatedly to sample from the distribution



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:

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

➤ Why? → Simulate with JUQCS

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

JUQCS

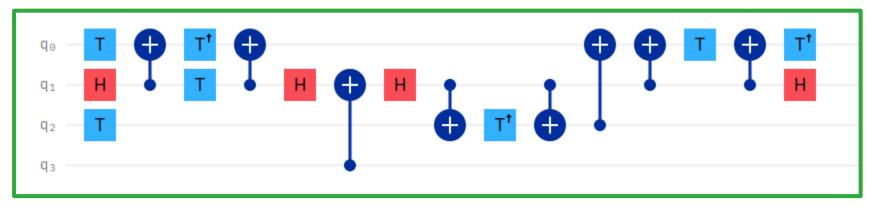
Jülich universal quantum computer simulator

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



Willsch et al. (2021

arXiv:2104.03293



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



> with huge vectors and huge² matrices



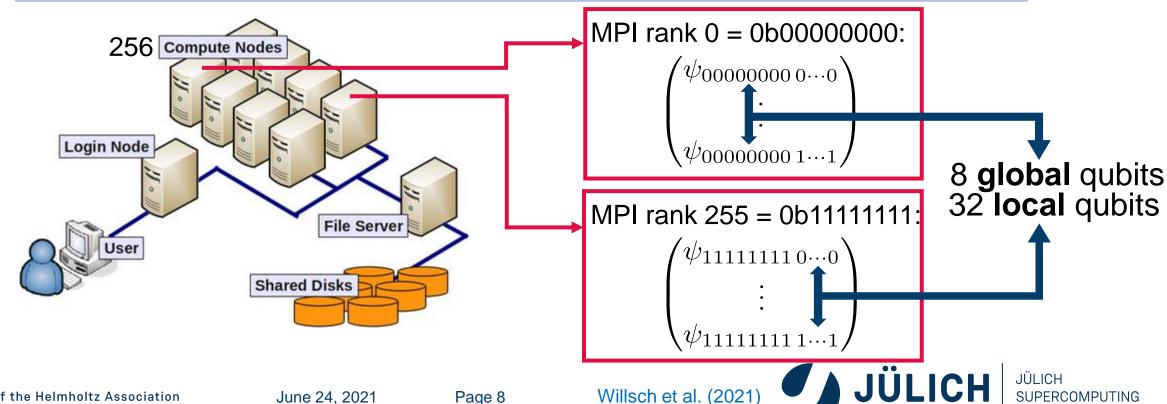
JUQCS

Distribution of the quantum state

How does the simulator manage all these complex numbers?

 \rightarrow Distribute quantum state $|\psi\rangle = (\psi_{\cdots q_2 q_1 q_0})$ over multiple compute nodes

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



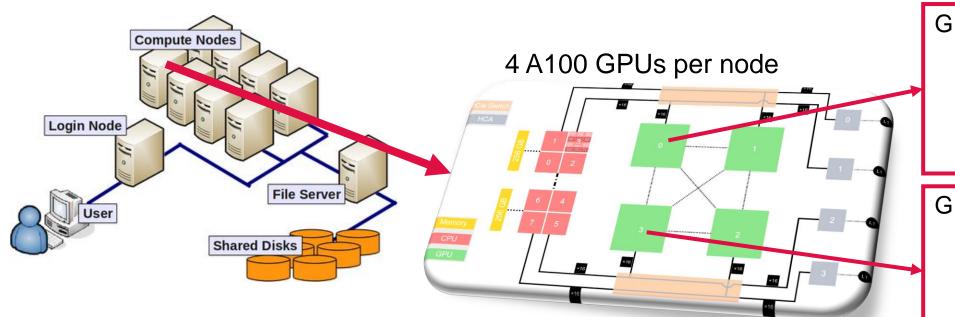
Willsch et al. (2021) arXiv:2104.03293



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



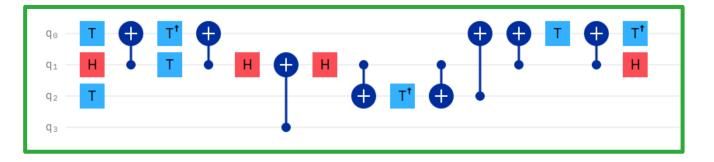
GPU rank 0 = 0b00000000000: $\psi_{000000000001\cdots 1}$.

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

> The MPI communication scheme is the same

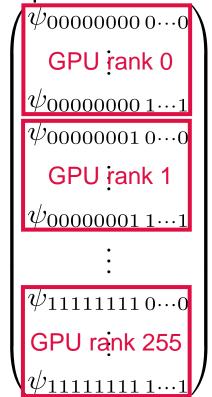
Willsch et al. (2021

CUDA implementation



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

Full quantum state:

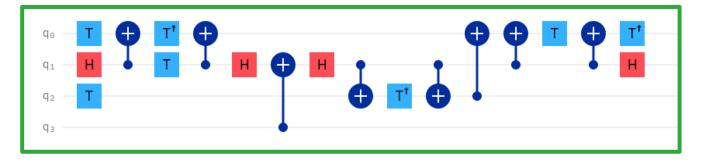


Quantum gate on **local** qubits:

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

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

CUDA implementation

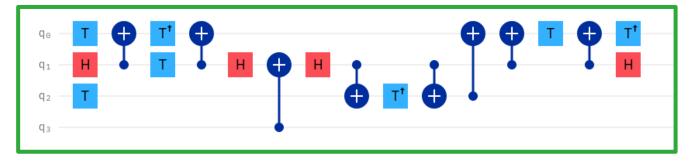


How to implement these matrix-vector multiplications in the most efficient way?

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

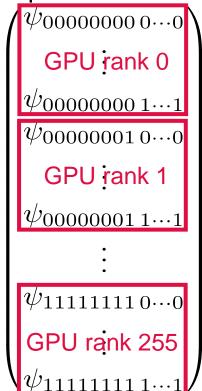
```
attributes(global) subroutine H operation GPU(nstates,psi R,psi I,ie,i1,c)
        USE cudafor
        implicit real(kind=8) (a-h,o-z)
        integer(kind=8),parameter :: one=1
        integer(kind=8), value :: nstates in
        integer(kind=8) :: m1,nm1,i,j,k
        integer(kind=4), value :: i0
        real(kind=8). 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)
        psi R(i)=(r0+r2)*c
        psi I(i)=(r1+r3)*c
        psi R(j)=(r0-r2)*c
 nd subroutine
```

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}

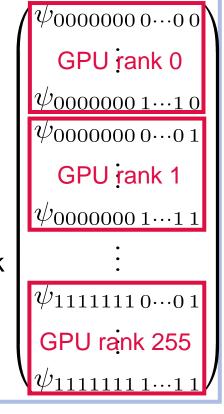
> Need to perform 2x2 updates of the form

H
$$\begin{pmatrix} \psi_{*********} & * \cdots * r \\ \psi_{*********} \end{pmatrix}$$

- > Problem: the numbers are on separate GPUs
- ➤ Naïve solution:
 - > Transfer $2^n/2 \ \psi's$ (8 TiB), perform \blacksquare , transfer back

Willsch et al. (2021)

- > 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

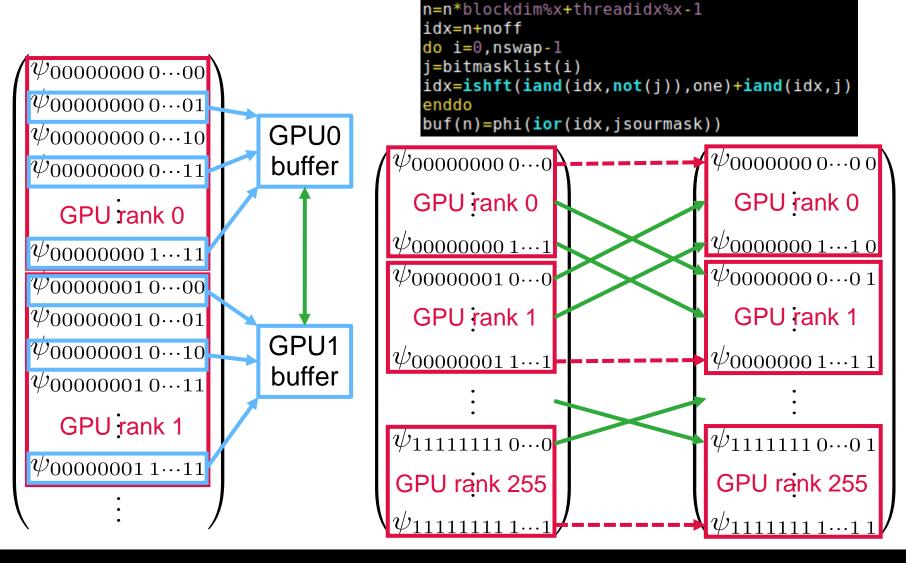
Before transfer:

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

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



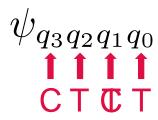
n=blockidx%x-1

```
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 MPIbutZps1<<<m,n>>> (nbits,nstatesGPU,m0,n0,nswap,isourmask,DD(0)%bitmasklist,DD(0)%phi_R,DD(0)%buf_R)
istat=cudaStreamSynchronize()
```

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 /



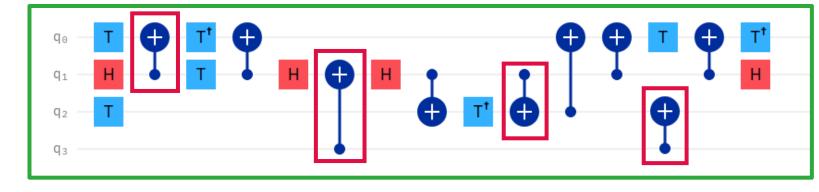


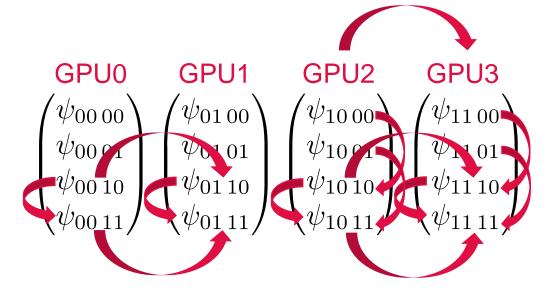
➤ "C global, T local": GPUs swap locally

➤ "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





JÜLICH

CENTRE

SUPERCOMPUTING

Running on JUWELS Booster

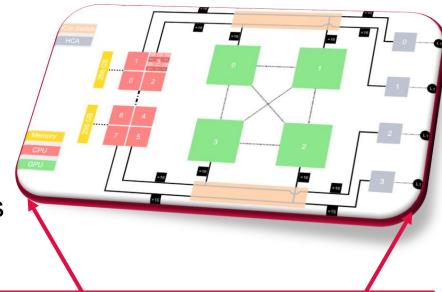


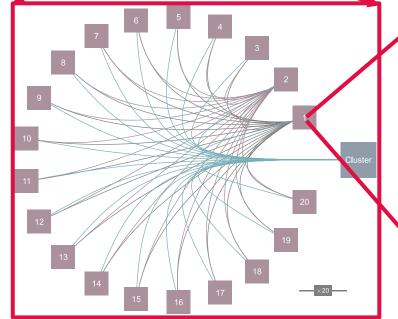
- Top500 Nov-2020:
- > #1 Europe
- > #7 World
- > #3 Green500

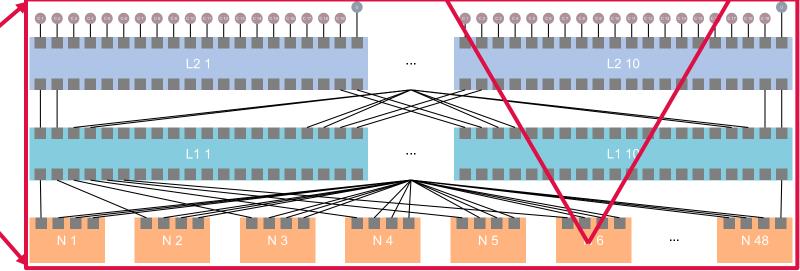


- ➤ 936 Compute nodes
- DragonFly+ Topology
- > 3744 NVIDIA A100 GPUs
- > 73 PFLOP/s

4 A100 GPUs per node







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

> 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:

Willsch et al. (2021

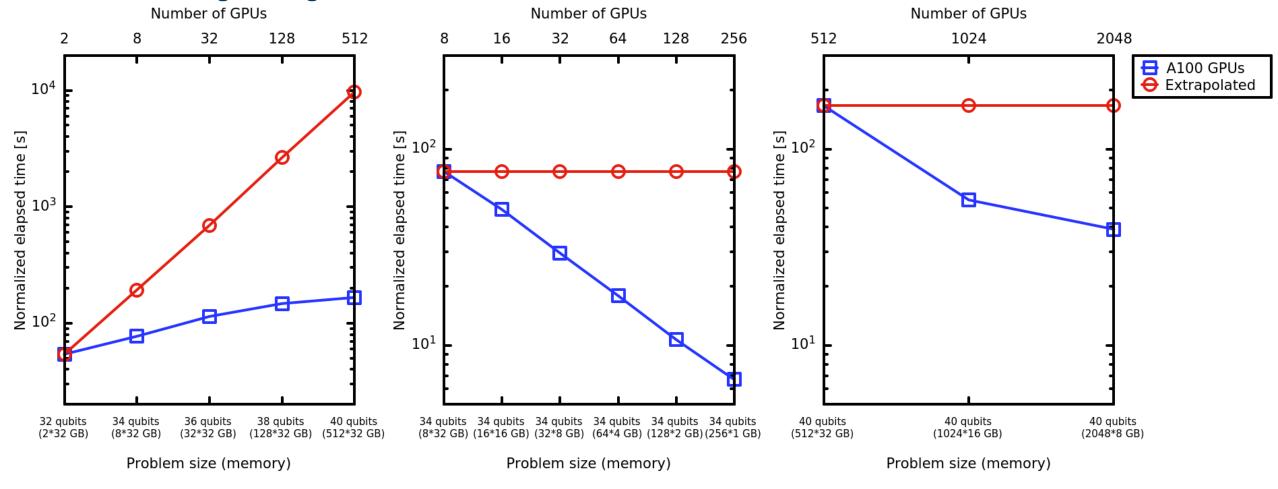
arXiv:2104.03293

Solve airplane scheduling problems using the Quantum Approximate Optimization Algorithm



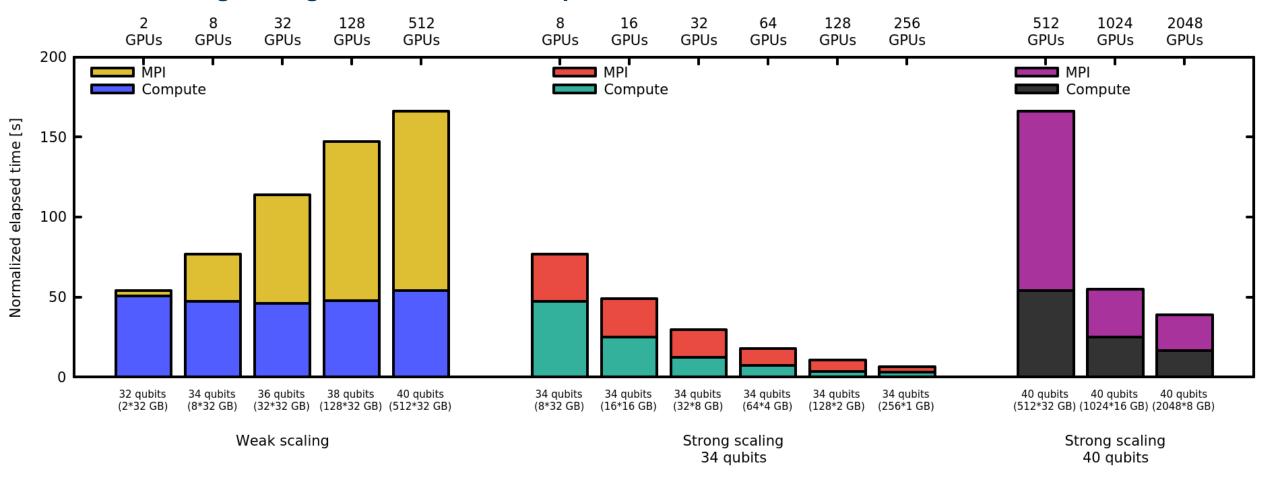
Willsch et al., <u>arXiv:2105.02208</u> (2021)

Weak and strong scaling results



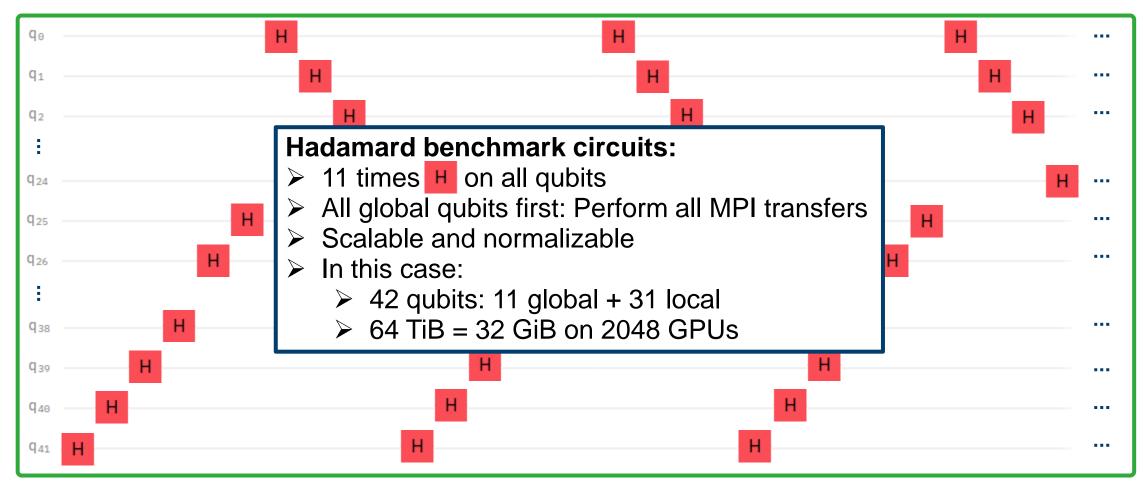
Willsch et al. (2021)

Weak and strong scaling results: MPI vs. Compute Time



Willsch et al. (2021)

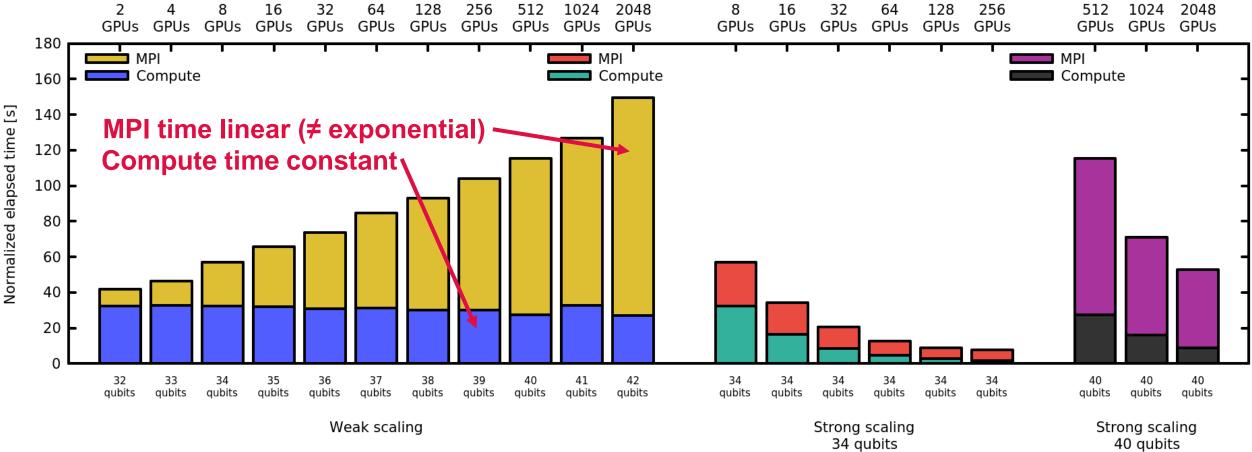
Weak and strong scaling results: Hadamard benchmark circuits



Benchmarking gate-based quantum computers: Michielsen et al., Comp. Phys. Commun. 220, 44 (2017)

Willsch et al. (2021

Weak and strong scaling results: Hadamard benchmark circuits

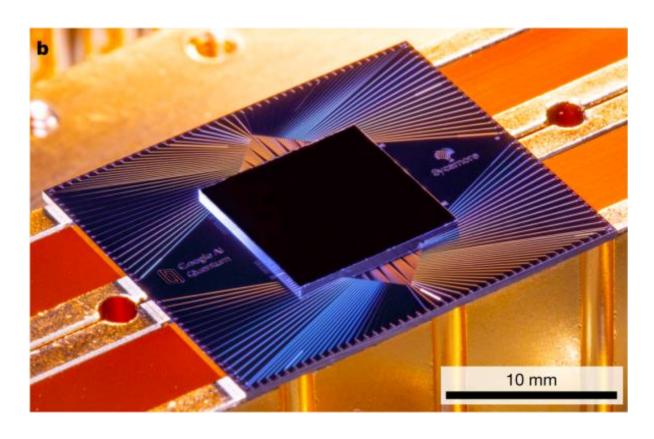


- > Speedup (compute time per node): JUWELS Cluster -> JUWELS GPUs (NVIDIA V100): 10
- ➤ Speedup (compute time per node): JUWELS GPUs → JUWELS Booster (NVIDIA A100): 2 3

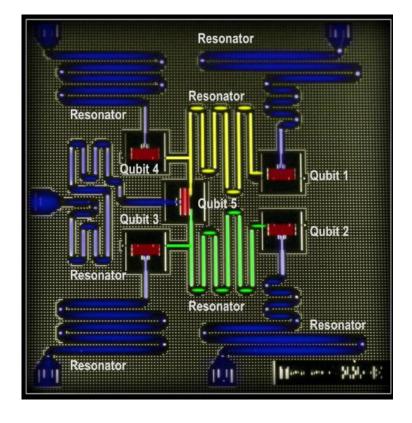
Willsch et al. (2021

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

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

- ➤ 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

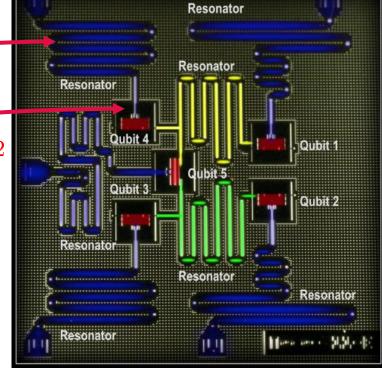
$$0 \le k_0, k_1 < 1000$$

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

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

Now:

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

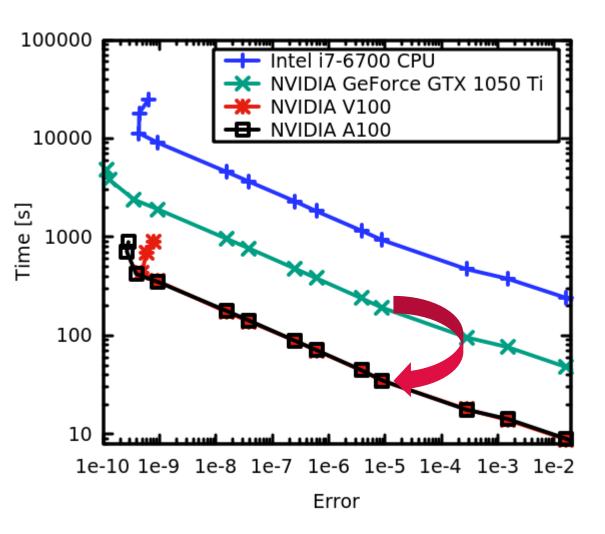


$$\tilde{c}_{k_0 m_0} = \cos(\tau \sqrt{k_0 + 1} (G\lambda_{m_0} + \varepsilon(\tilde{t})))$$

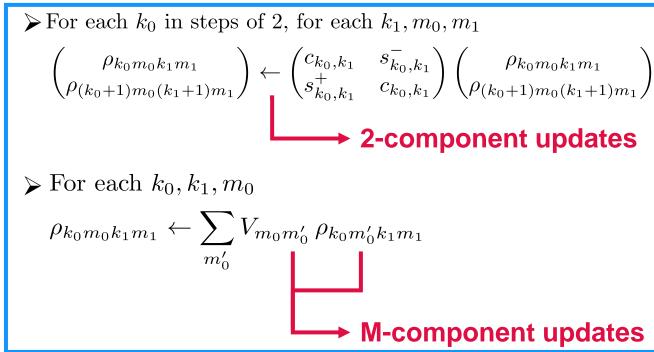
$$\tilde{s}_{k_0 m_0} = \sin(\tau \sqrt{k_0 + 1} (G\lambda_{m_0} + \varepsilon(\tilde{t})))$$

Simulating physical realizations of quantum computers on GPUs

CUDA+OpenACC MPI C++



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

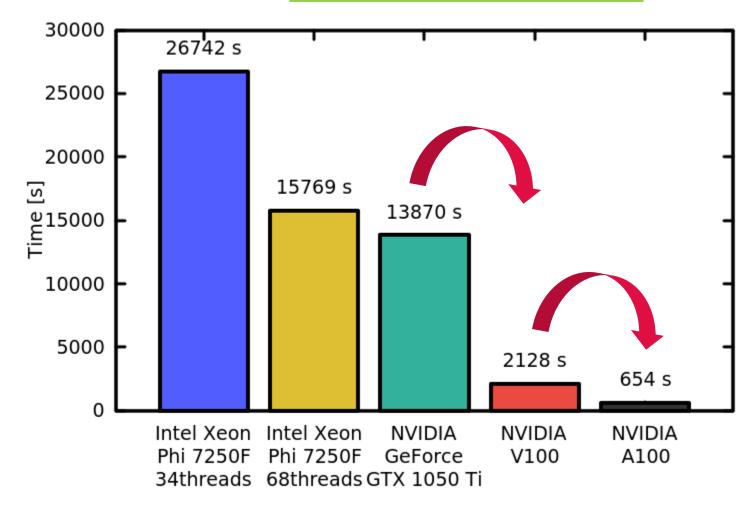
CUDA+OpenACC MPI C++

Realistic system with M=12 and K=400

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



Willsch et al. (2021)

arXiv:2104.03293

June 24, 2021

SUMMARY

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)
 - > Airplane scheduling problems: Willsch et al., https://arxiv.org/abs/2105.02208 (2021)
 - GPU-accelerated simulations of QA and the QAOA: Willsch et al., https://arxiv.org/abs/2104.03293 (2021)
 - > JUQMES: https://jugit.fz-juelich.de/gip/jugmes



Willsch et al. (2021

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

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

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

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

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

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

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

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

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

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

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

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

Willsch et al. (2021)

arXiv:2104.03293

→ 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

Willsch et al. (2021)

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

Yes, but:

> cuBLAS:

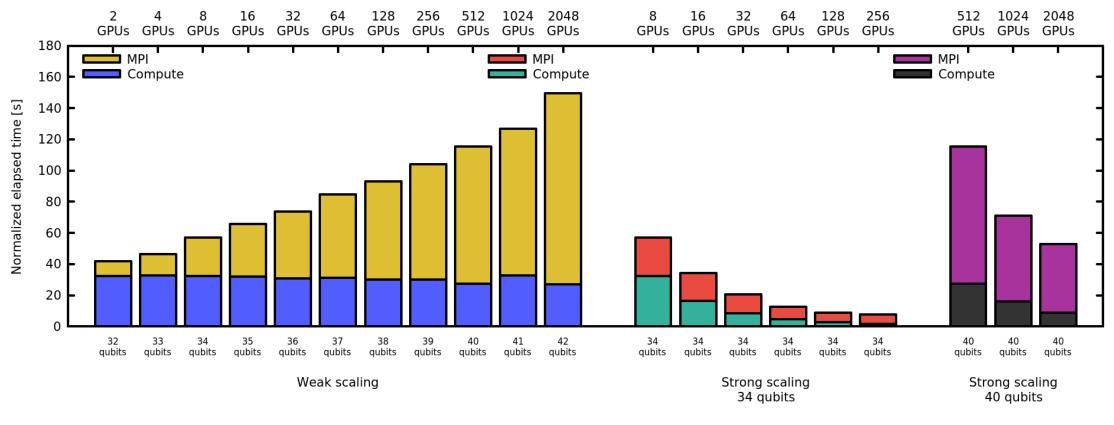
$$C = \alpha \operatorname{op}(A)\operatorname{op}(B) + \beta C$$

cublasGemmStridedBatchedEx:

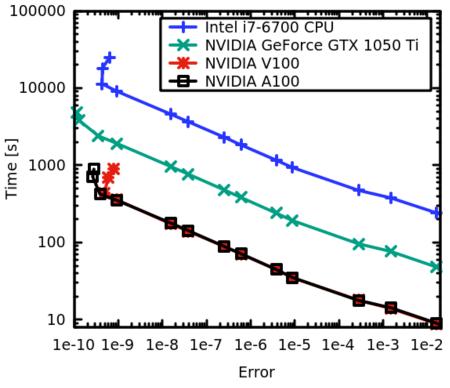
$$C + i * strideC = \alpha op(A + i * strideA)op(B + i * strideB) + \beta(C + i * strideC)$$

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

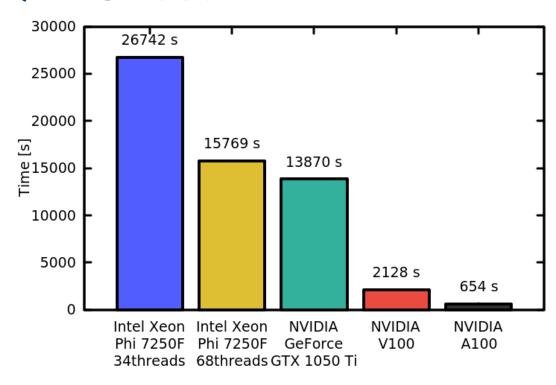
What is the bottleneck that limits the performance of JUQCS and are there ideas to improve it?



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



Small system with M=4 and K=100



Realistic system with M=12 and K=400



