A density-matrix renormalization group algorithm for simulating quantum circuits with a finite fidelity

Thomas Ayral Atos Quantum Laboratory

TA, Louvet, Zhou, Lambert, Stoudenmire, Waintal, 2207.05612



Closing the supremacy gap

Fall 2019:

Google claims quantum "supremacy"



+ others!





Today:

Even larger machines...



But no useful quantum advantage (yet).



What do we (really) need to beat to reach advantage?



Execute circuit (53 qubits, 230 two-qubit gates)



Return bitstrings x = 00101, ...



Arute et al '19



Execute circuit (53 qubits, 230 two-qubit gates)



Return bitstrings x = 00101, ...

Expectations





Arute et al '19

Atos

Execute circuit (53 qubits, 230 two-qubit gates)



Return bitstrings x = 00101, ...

Expectations





Arute et al '19



Execute circuit (53 qubits, 230 two-qubit gates)





Classical task: Draw *N*_{samples} with same XEB

6



What to expect?

The product law for the fidelity



Exponential decay!



What to expect?

The product law for the fidelity



1. Previous classical simulation strategies





Classical simulation of quantum circuits

From a tensor-network perspective

Goal: compute $P_U(x) = |\langle x | \Psi \rangle|^2 = |\langle x | U | 0 \rangle|^2$





Classical simulation of quantum circuits

From a tensor-network perspective

Goal: compute $P_U(x) = |\langle x | \Psi \rangle|^2 = |\langle x | U | 0 \rangle|^2$



 $s = \sum_{ijkl} A_{li} B_{ij} C_{jk} D_{kl}, \boldsymbol{O}(N^4)$

Clever contraction:

 $s = \sum_{l} \left[\sum_{k} \left\{ \sum_{j} \left(\sum_{i} A_{li} B_{ij} \right) C_{jk} \right\} D_{kl} \right], \boldsymbol{O}(N^3)$

Classical simulation of quantum circuits

From a tensor-network perspective

Goal: compute $P_U(x) = |\langle x | \Psi \rangle|^2 = |\langle x | U | 0 \rangle|^2$



Order matters:

- Naive contraction:
- $s = \sum_{ijkl} A_{li} B_{ij} C_{jk} D_{kl}, \boldsymbol{O}(\boldsymbol{N^4})$
- Clever contraction:

 $s = \sum_{l} \left[\sum_{k} \left\{ \sum_{j} \left(\sum_{i} A_{li} B_{ij} \right) C_{jk} \right\} D_{kl} \right], \boldsymbol{O}(N^3)$



Corresponding tensor network



• Find $\langle x|U|0 \rangle$: contract the tensor network.

Note: storage cost

- 3 qubits: 2x2x2 = 8
- *n* qubits: 2^{*n*}!



Three main classical simulation methods

... as three contraction strategies!

1. Schrödinger





Three main classical simulation methods

... as three contraction strategies!

Markov & Shi '08

1. Schrödinger



 C_1

2. "Tensor network"

- Find (close to) optimal contraction strategy (NP hard problem!)
- Contract (GPUs, TPUs...)

min(depth, width)

All amplitudes at once: "strong"



Three main classical simulation methods

... as three contraction strategies!

1. Schrödinger



All amplitudes at once: "strong"

2. "Tensor network"



 Find (close to) optimal contraction strategy (NP hard problem!)

2. Contract (GPUs, TPUs...)

3. Feynman (sum over paths)



 $\langle x|U|0\rangle = \sum_{\substack{a_1,a_2,\dots,c_1,c_2 \\ paths}} [\psi_0]_{a_1b_1c_1}[u_1]_{a_1a_2} \\ [\psi_0]_{a_2b_1,a_3b_2}[u_3]_{b_2c_1,b_3c_2} \\ \delta_{a_3x_1}\delta_{b_3x_2}\delta_{c_2x_3}$

depth



- **CPU**: exp(Treewidth)
- **Storage**: exp(Treewidth)

One amplitude: "closed"

- **CPU**: $\propto N_{\text{paths}} \sim \exp(N_{\text{gates}})$
- Storage: const.

One amplitude: "closed"

How did Google 'prove' supremacy?





Trick: simplified circuit

- easier contraction
- ... but same XEB!
 Test: full=simplified for smaller depths





How did Google 'prove' supremacy?

Arute et al '19



Need $P_U(x)!$

Trick: simplified circuit

- easier contraction
- ... but same XEB! Test: full=simplified for smaller depths



 a_1 a_2 a_3 a_1 b_2 b_3 a_2 c_2 a_3

For the true depth: simplified

STEP 2: Estimate classical time to get 0.2% XEB

• Essentially: Feynman approach.

Claim:

sum only fraction F of all paths => fidelity F

 $t_{\text{one path}} \times N_{\text{paths}} \times F$ = 10'000 years

... compare to: 200 secs to sample bitstrings

Goal: can we close this gap?



Previous attempts to catch up with Google

Pednault et al '19: Schrödinger, 2.5 days (est.)



Previous attempts to catch up with Google

Pednault et al '19: Schrödinger, 2.5 days (est.)

Gray & Kourtis '20, Huang et al '20: Index slicing: tensor network+ Feynman:



... 19.3 days (massively parallel)!

 $\langle x|U|0\rangle = \sum_{\mathbf{b}_{2}} T_{1}(x_{1}; \mathbf{b}_{2}) T_{2}(x_{2}, x_{3}; \mathbf{b}_{2})$

Previous attempts to catch up with Google

Pednault et al '19: Schrödinger, 2.5 days (est.)

Gray & Kourtis '20, Huang et al '20: Index slicing: tensor network+ Feynman:



... 19.3 days (massively parallel)!

See also Pan & Zhang '21 (5 days [GPUs]), Pan et al '21.

• ... down to 304 seconds (Liu et al '21)... on a very large HPC system (42 million cores)!

Main issue: exponential scaling! Very expensive to add a few qubits / gates...



Leveraging (low) entanglement





Beating the exponential with a finite fidelity Matrix Product States (MPS)

See e.g Schollwöck '11

• Previous attempts:

Surrender fidelity by summing fewer Feynman paths.

• New idea: use key quantum property: entanglement

Trivial case: Product states

 $[\psi_0]_{a_1b_1c_1} = [\psi_0^1]_{a_1} [\psi_0^2]_{b_1} [\psi_0^3]_{c_1}$



Beating the exponential with a finite fidelity

Matrix Product States (MPS)

• Example: an entangled state:

Previous attempts:

Surrender fidelity by summing fewer Feynman paths.

• New idea: use key quantum property: entanglement

Trivial case: Product states

 $[\psi_0]_{a_1b_1c_1} = [\psi_0^1]_{a_1} [\psi_0^2]_{b_1} [\psi_0^3]_{c_1}$

$[\psi_0]_{a_1b_1c_1} = \frac{1}{\sqrt{2}} [\psi_0^1]_{a_1} [\psi_0^2]_{b_1} [\psi_0^3]_{c_1} + \frac{1}{\sqrt{2}} [\chi_0^1]_{a_1} [\chi_0^2]_{b_1} [\chi_0^3]_{c_1}$

See e.q

Schollwöck '11



= "Matrix product state".

Beating the exponential with a finite fidelity

Matrix Product States (MPS)

Previous attempts:

Surrender fidelity by summing fewer Feynman paths.

• New idea: use key quantum property: entanglement

Trivial case: Product states

```
[\psi_0]_{a_1b_1c_1} = [\psi_0^1]_{a_1} [\psi_0^2]_{b_1} [\psi_0^3]_{c_1}
```

• Example: an entangled state:

 $[\psi_0]_{a_1b_1c_1} = \frac{1}{\sqrt{2}} [\psi_0^1]_{a_1} \, [\psi_0^2]_{b_1} \, [\psi_0^3]_{c_1} + \frac{1}{\sqrt{2}} \, [\chi_0^1]_{a_1} \, [\chi_0^2]_{b_1} \, [\chi_0^3]_{c_1}$



= "Matrix product state".



See e.g Schollwöck '11

Zhou, Stoudenmire, Waintal PRX 20

Algorithm to compute $\langle x|U|0\rangle$: Vidal '04





Algorithm to compute $\langle x|U|0\rangle$: Vidal '04





Algorithm to compute $\langle x|U|0\rangle$: Vidal '04





Zhou, Stoudenmire, Waintal PRX 20





Vidal '04

Zhou, Stoudenmire, Waintal PRX 20

SVD compression J1 SVD compression f_2 \approx

• Final fidelity: $F = f_1 f_2$

Algorithm to compute $\langle x|U|0\rangle$:

• Works... but not enough to reproduce Google



Algorithm to compute $\langle x|U|0\rangle$: Vidal '04 SVD compression SVD compression f_2 \approx

Improvement: Schrödinger + MPS

• "Group tensors together"



Useful for 2D qubit grids:

- Vertical gates: "exact"
- Horizontal gates: "compressed"



• Final fidelity: $F = f_1 f_2$

• Works... but not enough to reproduce Google

Improves fidelity... but still not enough!



This work: a triply hybrid strategy

MPS + Schrödinger + tensor networks via a Density Matrix Renormalization Group method

Previous approach: apply 1 gate and compress

Here: apply several layers of gates, ... and find "optimal" MPS:





This work: a triply hybrid strategy

MPS + Schrödinger + tensor networks via a Density Matrix Renormalization Group method

Previous approach: apply 1 gate and compress

Here: apply several layers of gates, ... and find "optimal" MPS: How to find optimal MPS?

- DMRG: find MPS with maximal overlap
- Tensor-by-tensor optimization: n_s "sweeps"

TA, Louvet, Zhou, Lambert, Stoudenmire, Waintal, 2207.05612





Contracting the tensor network

Mix of automatic contraction with "manual" guidance:





3. Results





A different relation between F and XEB

Reminder: in noisy random circuits:

 $F \approx XEB$

In our method:

XEB $\approx \sqrt{F}$

• Strange?



A different relation between F and XEB

Reminder: in noisy random circuits:

 $F \approx XEB$

In our method:

XEB $\approx \sqrt{F}$

• Strange?

In the chaotic limit (infinite depth):

$$F = \frac{1}{2^n}$$
$$XEB = \frac{1}{2^{n/2}}$$



A different relation between F and XEB

Atos QLM





Reminder: in noisy random circuits:

 $F \approx XEB$

In our method:

XEB $\approx \sqrt{F}$

• Strange?

In the chaotic limit (infinite depth):

$$F = \frac{1}{2^n}$$
$$XEB = \frac{1}{2^{n/2}}$$

Atos

Closing the supremacy gap

Error per gate: $\varepsilon = 1 - \tilde{F}^{1/N_{2g}}$





Closing the supremacy gap

Error per gate: $\varepsilon = 1 - \tilde{F}^{1/N_{2g}}$



Closed (one amplitude) vs open (all):

- Open: most powerful. Strong simulation.
- Closed: can contract from "both sides". Lower compression losses!



"Closed" easily usable for "weak" simulation: Bravyi et al '21.

Here (near-chaotic distribution): Could use Metropolis-Hastings: 70% acceptance rate.



The influence of the type of circuit



The influence of the type of circuit



A scalable method: what happens when increasing the qubit count?

Fixed depth *D*:

• Error per gate increases... then stagnates:



But more gates, XEB decreases...: must increase N_{samples} to reduce variance.

Keep nD fixed (fixed XEB: fixed experimental time!):

• better and better error per gate!



A scalable method: what happens when increasing the qubit count?

Fixed depth *D*:

• Error per gate increases... then stagnates:



But more gates, XEB decreases...: must increase $N_{
m samples}$ to reduce variance.

Keep nD fixed (fixed XEB: fixed experimental time!):

• better and better error per gate!

How to understand the stagnation?

Can compute "optimal" error rate after SVD compression:

$$\varepsilon_{\rm opt} = \frac{1}{D} \left(\log 2 - \frac{\log 4\chi}{2N} \right)$$



when ε_{SVD} reaches ε_{opt} [~chaotic limit], actual error rate deviates from ε_{SVD}



Conclusions

Available on Atos QLM ("QPEG")



Google's game can be won by a classical simulator

- Key: leverage 3 different simulations methods:
 - Schrödinger, tensor-network, Matrix Product State
- Very basic implementation (no sophisticated contraction)

Beware of "large Hilbert space fallacy": Structure matters!

- MPS: capture true difficult quantity: entanglement
- This is the target for quantum advantage!

Qubit quality is crucial! Recent progress in QEC @ETH, Google, Quantinuum...

XEB flaws: not scalable, not 'useful' task

• Q-score (Martiel, TA, Allouche '20): use QAOA MaxCut as benchmark

