# Evaluating the performance of quantum processing units (QPUs) at large width and depth
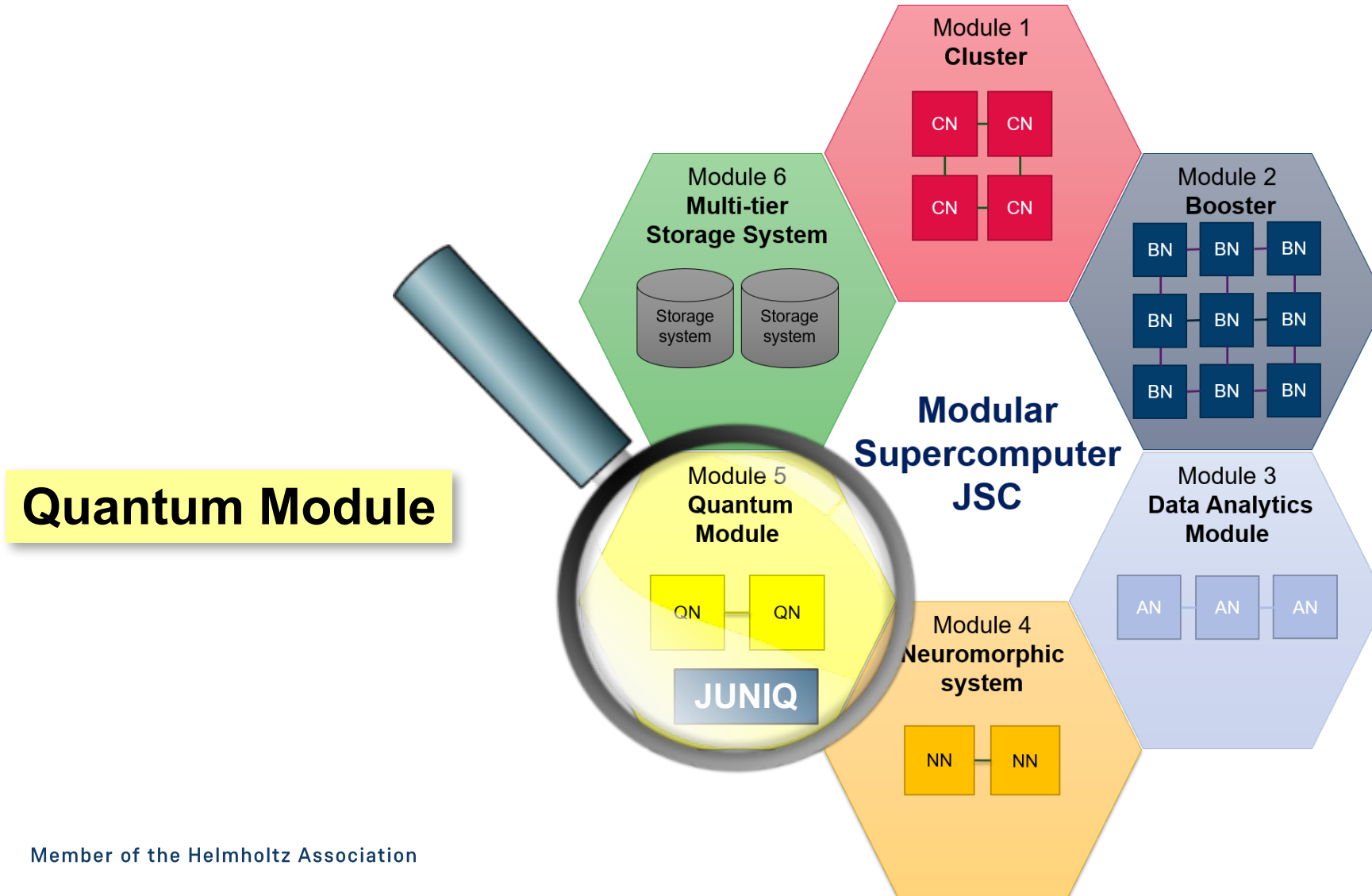
TQCI seminar dedicated to benchmarks for quantum computers

June 16, 2025 I Alejandro Montanez-Barrera | Forschungszentrum Jülich - Jülich Supercomputing Center
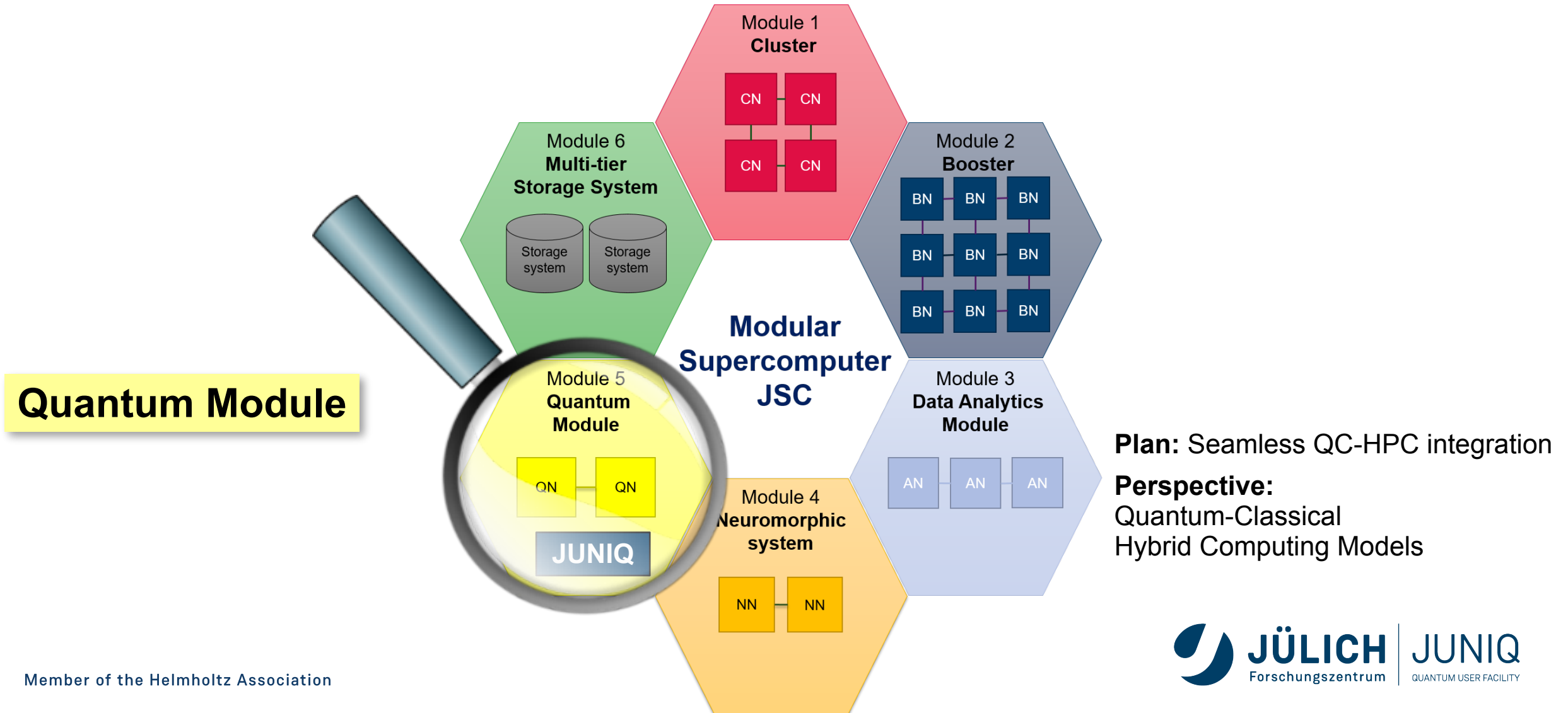
https://arxiv.org/abs/2502.06471

Member of the Helmholtz Association

JÜLICH
Forschungszentrum

# JUNIQ - Jülich UNified Infrastructure for Quantum computing

## A European quantum computer user facility at the Jülich Supercomputing Centre



**Quantum Module**

Module 1
**Cluster**
CN CN
CN CN

Module 6
**Multi-tier Storage System**
Storage system  Storage system

Module 2
**Booster**
BN BN BN
BN BN BN
BN BN BN

**Modular Supercomputer JSC**

Module 5
**Quantum Module**
QN QN
JUNIQ

Module 3
**Data Analytics Module**
AN AN AN

Module 4
**Neuromorphic system**
NN NN

**Plan:** Seamless QC-HPC integration

**Perspective:** Quantum-Classical Hybrid Computing Models

JÜLICH Forschungszentrum | JUNIQ QUANTUM USER FACILITY

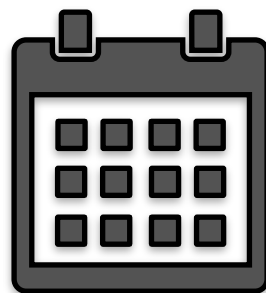# JUNIQ - Jülich UNified Infrastructure for Quantum computing



1. QC user facility for science and industry
2. Installation, operation and provision of QCs
3. Unified portal for access to QC emulators and to QC devices at different levels of technological maturity.
4. Development of algorithms and prototype applications
5. Services, training and user support
6. Modular quantum-HPC hybrid computing
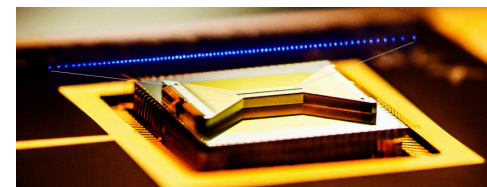
# Why do we do quantum benchmarking?

Quantum benchmarking is essential because it provides **quantitative, standardized ways to evaluate and compare quantum devices** and **algorithms**. In short, it tells us **how good a quantum computer really is** and whether it's improving.

**Measure Device Performance**

**Track Progress Over Time**

**Compare Across Platforms**

JÜLICH
Forschungszentrum

# How to measure the progress of quantum processing units?

Randomized benchmarking is proposed by Dankert et al

2006

https://arxiv.org/abs/quant-ph/0606161

JÜLICH
Forschungszentrum

# How to measure the progress of quantum processing units?

aws

Average 2Q fidelity (%) Info

98.720

**Average 2Q fidelity (%)**

Average gate fidelities for a maximally entangling two-qubit native gate. It is regularly measured by IonQ on the device using Direct Randomized Benchmarking

Randomized benchmarking is proposed by Dankert et al

**2006**

https://arxiv.org/abs/quant-ph/0606161

JÜLICH
Forschungszentrum

# How to measure the progress of quantum processing units?

aws

Average 2Q fidelity (%) **Info**

98.720

**Average 2Q fidelity (%)**

Average gate fidelities for a maximally entangling two-qubit native gate. It is regularly measured by IonQ on the device using Direct Randomized Benchmarking

CZ fidelity (%)

IQM 99.232

CZ gate fidelities for the qubit pair measured by interleaved randomized benchmarking on 2 or 3 pairs simultaneously.

Randomized benchmarking is proposed by Dankert et al

**2006**

https://arxiv.org/abs/quant-ph/0606161

JÜLICH
Forschungszentrum

# How to measure the progress of quantum processing units?

aws

Average 2Q fidelity (%) [Info]

98.720

**Average 2Q fidelity (%)**

Average gate fidelities for a maximally entangling two-qubit native gate. It is regularly measured by IonQ on the device using Direct Randomized Benchmarking

CZ fidelity (%)

IQM 99.232

CZ gate fidelities for the qubit pair measured by interleaved randomized benchmarking on 2 or 3 pairs simultaneously.

IBM Quantum **Platform**

Median CZ error:

2.696e-3      99.73 %

2Q error (best): The lowest two-qubit error rate from all edges measured by isolated randomized benchmarking.

Randomized benchmarking is proposed by Dankert et al

**2006**

https://arxiv.org/abs/quant-ph/0606161

JÜLICH
Forschungszentrum

# How to measure the progress of quantum processing units?

aws

**Average 2Q fidelity (%)** Info

98.720

**Average 2Q fidelity (%)**

Average gate fidelities for a maximally entangling two-qubit native gate. It is regularly measured by IonQ on the device using Direct Randomized Benchmarking

CZ fidelity (%)

IQM 99.232

CZ gate fidelities for the qubit pair measured by interleaved randomized benchmarking on 2 or 3 pairs simultaneously.

QUANTINUUM SYSTEMS

**2-qubit gate fidelity**

The 2-qubit gate fidelity is measured with 2-qubit RB

$$H2 - 1 : 99.89\,\%$$

IBM Quantum **Platform**

Median CZ error:

2.696e-3       99.73 %

2Q error (best): The lowest two-qubit error rate from all edges measured by isolated randomized benchmarking.

Randomized benchmarking is proposed by Dankert et al

**2006**

https://arxiv.org/abs/quant-ph/0606161

JÜLICH
Forschungszentrum

# How to measure the progress of quantum processing units?

**aws**

Average 2Q fidelity (%) [Info]

98.720

**Average 2Q fidelity (%)**

Average gate fidelities for a maximally entangling two-qubit native gate. It is regularly measured by IonQ on the device using Direct Randomized Benchmarking

CZ fidelity (%)

**IQM** 99.232

CZ gate fidelities for the qubit pair measured by interleaved randomized benchmarking on 2 or 3 pairs simultaneously.

**QUANTINUUM SYSTEMS**

**2-qubit gate fidelity**

The 2-qubit gate fidelity is measured with 2-qubit RB

$$H2 - 1 : 99.89\,\%$$

**rigetti**

ISWAP gate fidelity (%)

| Average | Median |
|---------|--------|
| 94.431 | 98.700 |

**IBM Quantum Platform**

Median CZ error:

2.696e-3    99.73 %

2Q error (best): The lowest two-qubit error rate from all edges measured by isolated randomized benchmarking.
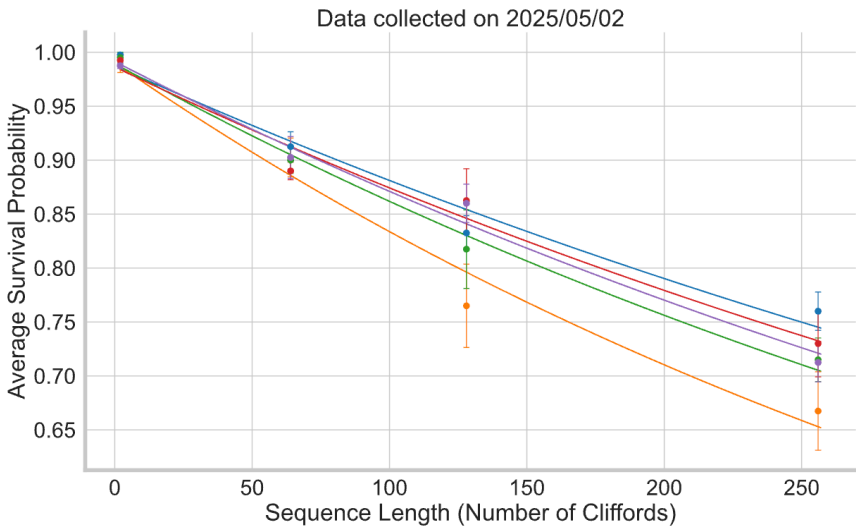
Randomized benchmarking is proposed by Dankert et al

**2006**

https://arxiv.org/
abs/quant-ph/
0606161

**JÜLICH** Forschungszentrum

# How to measure the progress of quantum processing units?

**aws**

Average 2Q fidelity (%) [Info]

98.720

**Average 2Q fidelity (%)**

Average gate fidelities for a maximally entangling two-qubit native gate. It is regularly measured by IonQ on the device using Direct Randomized Benchmarking

CZ fidelity (%)

IQM 99.232

CZ gate fidelities for the qubit pair measured by interleaved randomized benchmarking on 2 or 3 pairs simultaneously.

IBM Quantum **Platform**

Median CZ error:

2.696e-3    99.73 %

2Q error (best): The lowest two-qubit error rate from all edges measured by isolated randomized benchmarking.

QUANTINUUM **SYSTEMS**

**2-qubit gate fidelity**

The 2-qubit gate fidelity is measured with 2-qubit RB

$$H2 - 1 : 99.89\,\%$$

**rigetti**

ISWAP gate fidelity (%)

| Average | Median |
|---------|--------|
| 94.431  | 98.700 |

Data collected on 2025/05/02



Randomized benchmarking is proposed by Dankert et al

**2006**

https://arxiv.org/abs/quant-ph/0606161

JÜLICH
Forschungszentrum

# How to measure the progress of quantum processing units?

Randomized benchmarking is proposed by Dankert et al
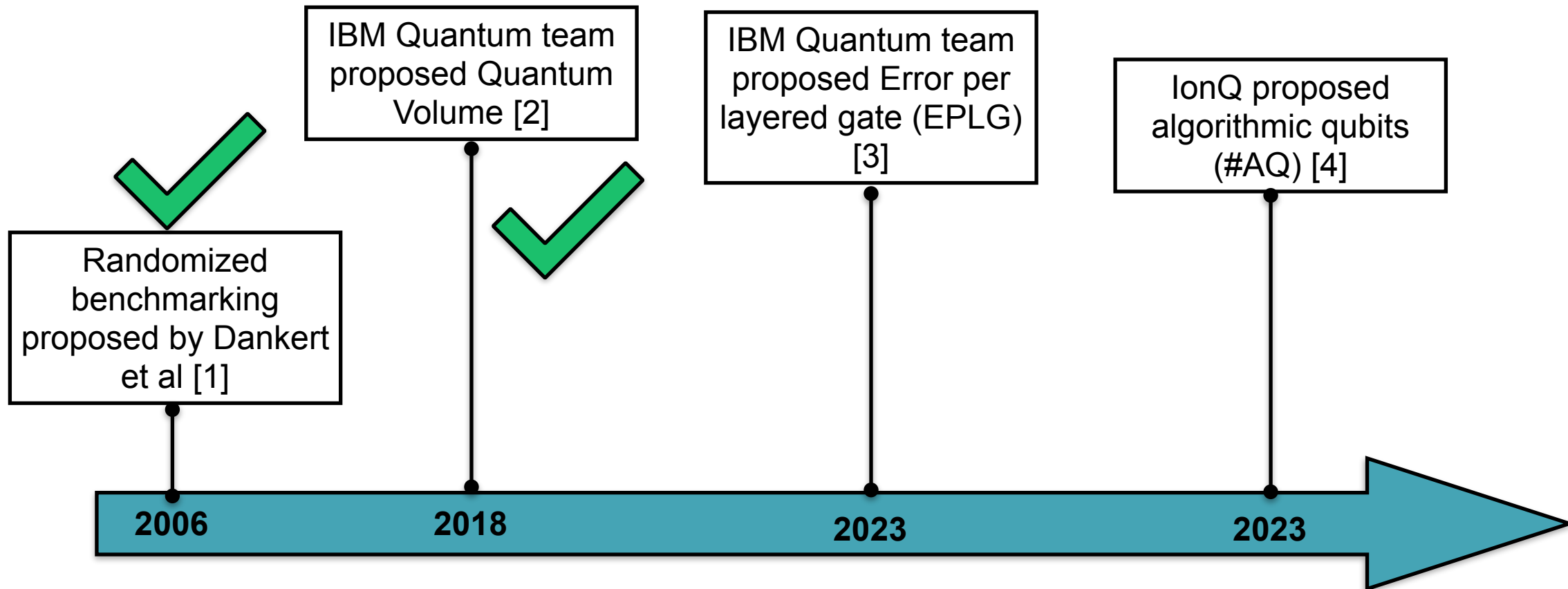
2006

https://arxiv.org/abs/quant-ph/0606161

JÜLICH
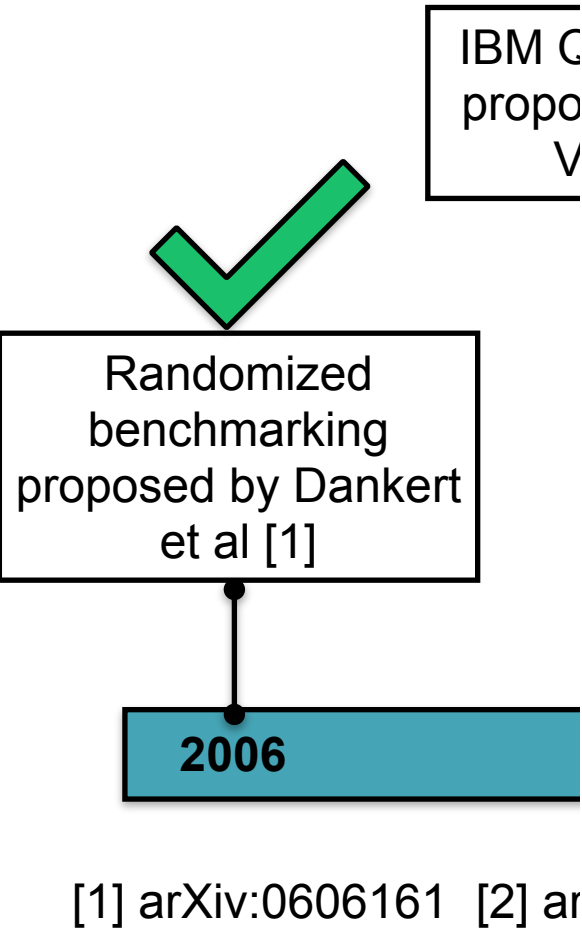Forschungszentrum

# How to measure the progress of quantum processing units?



IBM Quantum team proposed Quantum Volume [2]

Randomized benchmarking proposed by Dankert et al [1]

2006

2018

[1] https://arxiv.org/abs/quant-ph/0606161

[2] https://arxiv.org/pdf/1811.12926

JÜLICH
Forschungszentrum

# How to measure the progress of quantum processing units?

IBM Quantum team proposed Quantum Volume [2]

Randomized benchmarking proposed by Dankert et al [1]

Random permutation

two-qubit gate

$d \approx N_q$

$N_q$

$\pi$

$$\log_2 V_Q = \underset{N_q}{\operatorname{argmax}} \min(N_q, d(N_q))$$

**2006**        **2018**

[1] https://arxiv.org/abs/quant-ph/0606161

[2] https://arxiv.org/pdf/1811.12926

JÜLICH
Forschungszentrum

# How to measure the progress of quantum processing units?

IBM Quantum team proposed Quantum Volume [2]

Random permutation

two-qubit gate

$d \approx N_q$

Randomized benchmarking proposed by Dankert et al [1]

$N_q$

$\pi$

$$\log_2 V_Q = \underset{N_q}{\text{argmax}} \min(N_q, d(N_q))$$

**2006**　　　　**2018**

[1] https://arxiv.org/abs/quant-ph/0606161

[2] https://arxiv.org/pdf/1811.12926

$h_U \approx 0.85$ ideal

$h_U > 2/3$ pass test

$h_U = 0.5$ fully mixed

JÜLICH
Forschungszentrum

# How to measure the progress of quantum processing units?

Random permutation

two-qubit gate

$d \approx N_q$

IBM Quantum team proposed Quantum Volume [2]

$N_q$

$\pi$

Randomized benchmarking proposed by Dankert et al [1]

$$\log_2 V_Q = \underset{N_q}{\mathrm{argmax}} \min(N_q, d(N_q))$$

**2006**　　　　**2018**

[1] https://arxiv.org/abs/quant-ph/0606161

[2] https://arxiv.org/pdf/1811.12926

$h_U \approx 0.85$ ideal

$h_U > 2/3$ pass test

$h_U = 0.5$ fully mixed

$N_q$

$V_Q$

$d$

JÜLICH
Forschungszentrum

# How to measure the progress of quantum processing units?

**IBM Quantum team proposed Quantum Volume [2]**

**Randomized benchmarking proposed by Dankert et al [1]**

Random permutation

two-qubit gate

$d \approx N_q$

$N_q$

$\pi$

$$\log_2 V_Q = \underset{N_q}{\operatorname{argmax}} \min(N_q, d(N_q))$$

**Define a clear set of rules!**

"The transpiler is free to use all available tricks and hardware resources to implement U (e.g., taking great computational effort in finding an optimized U, using extra qubits for gate teleportation or temporary storage, etc.). It may optimize over qubit placements by choosing the best region of the device. If it is practical to calibrate a very large gate set, and it happens to include an accurate implementation of U, the transpiler is free to use it."[2]

**2006**          **2018**

[1] https://arxiv.org/abs/quant-ph/0606161

[2] https://arxiv.org/pdf/1811.12926

$h_U \approx 0.85$ ideal

$h_U > 2/3$ pass test

$h_U = 0.5$ fully mixed

$N_q$

$V_Q$

$d$

**JÜLICH**
Forschungszentrum

# How to measure the progress of quantum processing units?

**IBM Quantum team proposed Quantum Volume [2]**

**Randomized benchmarking proposed by Dankert et al [1]**

Random permutation

two-qubit gate

$d \approx N_q$

$N_q$

$\pi$

$$\log_2 V_Q = \underset{N_q}{\mathrm{argmax}} \min(N_q, d(N_q))$$

**2006**

**2018**

[1] https://arxiv.org/abs/quant-ph/0606161

[2] https://arxiv.org/pdf/1811.12926

**Define a clear set of rules!**

"The transpiler is free to use all available tricks and hardware resources to implement U (e.g., taking great computational effort in finding an optimized U, using extra qubits for gate teleportation or temporary storage, etc.). It may optimize over qubit placements by choosing the best region of the device. If it is practical to calibrate a very large gate set, and it happens to include an accurate implementation of U, the transpiler is free to use it."[2]

**Limitation**

"To determine if a given output is heavy, we compute $H_U$ directly from U using a method that scales exponentially with Nq".

$h_U \approx 0.85$ ideal

$h_U > 2/3$ pass test

$h_U = 0.5$ fully mixed

$N_q$

$V_Q$

$d$

**JÜLICH** Forschungszentrum

# How to measure the progress of quantum processing units?

IBM Quantum team proposed Quantum Volume [2]

Randomized benchmarking proposed by Dankert et al [1]



**2006**

**2018**

[1] https://arxiv.org/abs/quant-ph/0606161

[2] https://arxiv.org/pdf/1811.12926

JÜLICH
Forschungszentrum

# How to measure the progress of quantum processing units?

Randomized benchmarking proposed by Dankert et al [1]

IBM Quantum team proposed Quantum Volume [2]

IBM Quantum team proposed Error per layered gate (EPLG) [3]

**2006**   **2018**   **2023**

[1] arXiv:0606161   [2] arXiv:1811.12926   [3] arXiv:2311.05933

Error per layered gate

# How to measure the progress of quantum processing units?



Randomized benchmarking proposed by Dankert et al [1]

IBM Quantum team proposed Quantum Volume [2]

IBM Quantum team proposed Error per layered gate (EPLG) [3]

IonQ proposed algorithmic qubits (#AQ) [4]

2006

2018

2023

2023

[1] arXiv:0606161  [2] arXiv:1811.12926     [3] arXiv:2311.05933     [4] arXiv:2308.05071

JÜLICH
Forschungszentrum

# How to measure the progress of quantum processing units?



[1] arXiv:0606161  [2] ar...

[4] arXiv:2308.05071

# How to measure the progress of quantum processing units?

IBM ...
prop...
V...

Randomized
benchmarking
proposed by Dankert
et al [1]

**2006**

[1] arXiv:0606161  [2] ar...

IonQ proposed
algorithmic qubits
(#AQ) [4]

**Debunking algorithmic qubits**

QUANTINUUM

**2023**

[4] arXiv:2308.05071

#AQ Benchmark on IonQ Aria (Merged)
Nov 13, 2022

Hamiltonial Simulation
Phase Estimation
Quantum Fourier Tansform (1)
VQE Simulation (1)
Amplitude Estimation
Monte Carlo Sampling (2)

Circuit Width

Avg Result

AQ=25.0

Number of 2Q gates

JÜLICH
Forschungszentrum

# The Linear Ramp Quantum Approximate Optimization Algorithm (LR-QAOA)



(a) Graph topologies
(b) QAOA algorithm
(c) Linear ramp protocol for QAOA
(d) Expected performance with noise

We applied this benchmark methodology to 24 different QPUs from 6 different vendors, IQM, IBM, Rigetti, IonQ, Quantinuum, and OriginQ using 5 to 156 qubits and up to p=10,000.

JÜLICH
Forschungszentrum

# The problem behind LR-QAOA

$$x* = 00101$$



The weighted maxcut (WMC) problem involves determining the partition of the vertices in an undirected graph so that the total weight of the edges between the two sets is maximized.

**Cost function**

$$C(x) = \sum_{(i,j)\in E} w_{ij}(x_i + x_j - 2x_i x_j)$$

**Approximation ratio**

$$r = \frac{\sum_{k=1}^{n} C(x^k)/n}{C(x*)}$$

$x_k$    sample solution

$x*$    optimal solution

$n$    Samples

JÜLICH
Forschungszentrum

# The problem behind LR-QAOA

$$x^* = 00101$$



The weighted maxcut (WMC) problem involves determining the partition of the vertices in an undirected graph so that the total weight of the edges between the two sets is maximized.

**Cost function**

$$C(x) = \sum_{(i,j)\in E} w_{ij}(x_i + x_j - 2x_i x_j)$$

**Approximation ratio**

$$r = \frac{\sum_{k=1}^{n} C(x^k)/n}{C(x^*)}$$

$x_k$     sample solution

$x^*$     optimal solution

$n$     Samples

Performance metric

JÜLICH
Forschungszentrum

# Tracking the evolution of real QPUs

# Distinguishing successful results

To certify if the result of a QPU is still meaningful, we compare the approximation ratio for the LR-QAOA WMC problem given by the samples of the QPU to those coming from a random sampler.



(a) H2-1 50-qubit, 50 samples, and p=4.

# LR-QAOA on a 1D-Chain graph



Performance on a 100-qubit 1D chain experiment. IBM QPUs (EPLG)

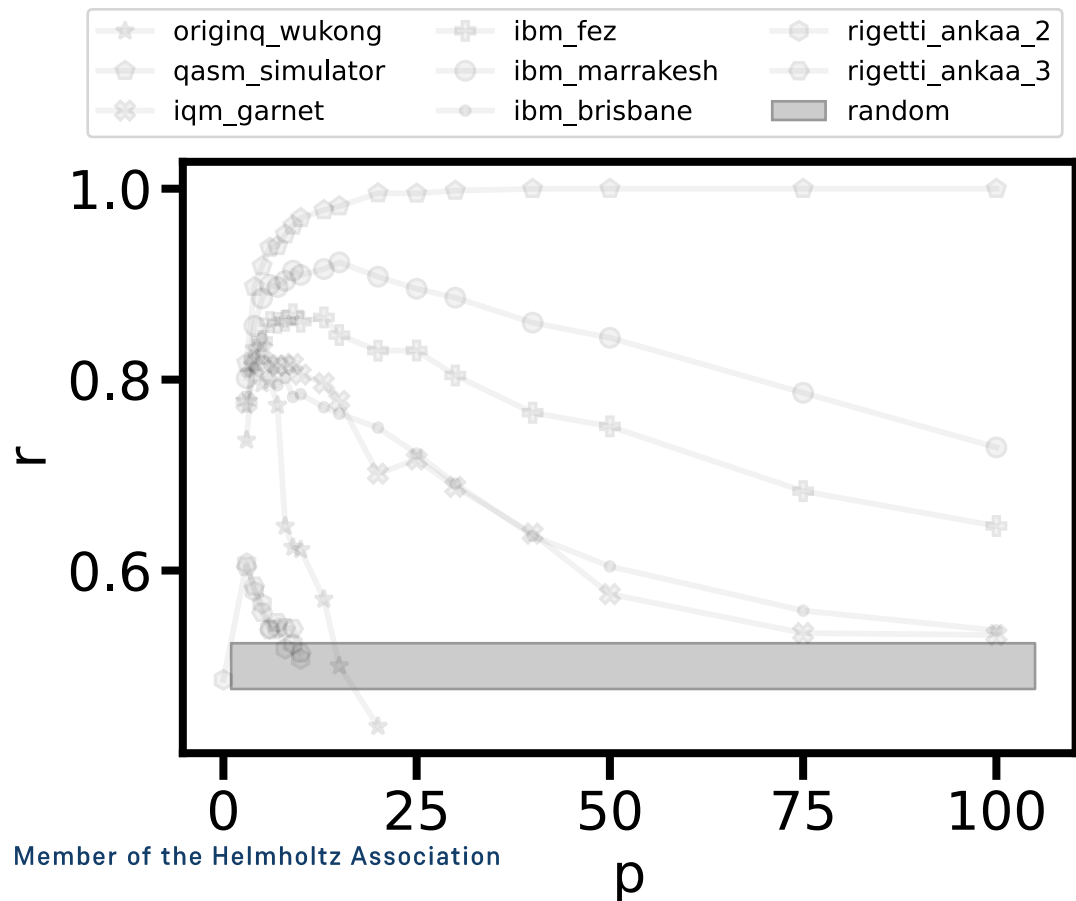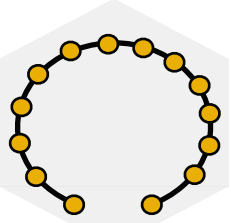Performance on the best 5-qubit 1D chain experiment, different QPUs

Legend (100-qubit): ibm_marrakesh (0.4), ibm_fez (0.8), ibm_torino-v1 (1.1), ibm_torino-v0 (0.8), ibm_brisbane (1.9), ibm_sherbrooke (1.7), ibm_kyiv (2.1), ibm_nazca (3.2), ibm_kyoto (3.6), ibm_osaka (2.8), ibm_brussels (2.2), ibm_strasbourg (5.4), random

Legend (5-qubit): originq_wukong, qasm_simulator, iqm_garnet, ibm_fez, ibm_marrakesh, ibm_brisbane, rigetti_ankaa_2, rigetti_ankaa_3, random
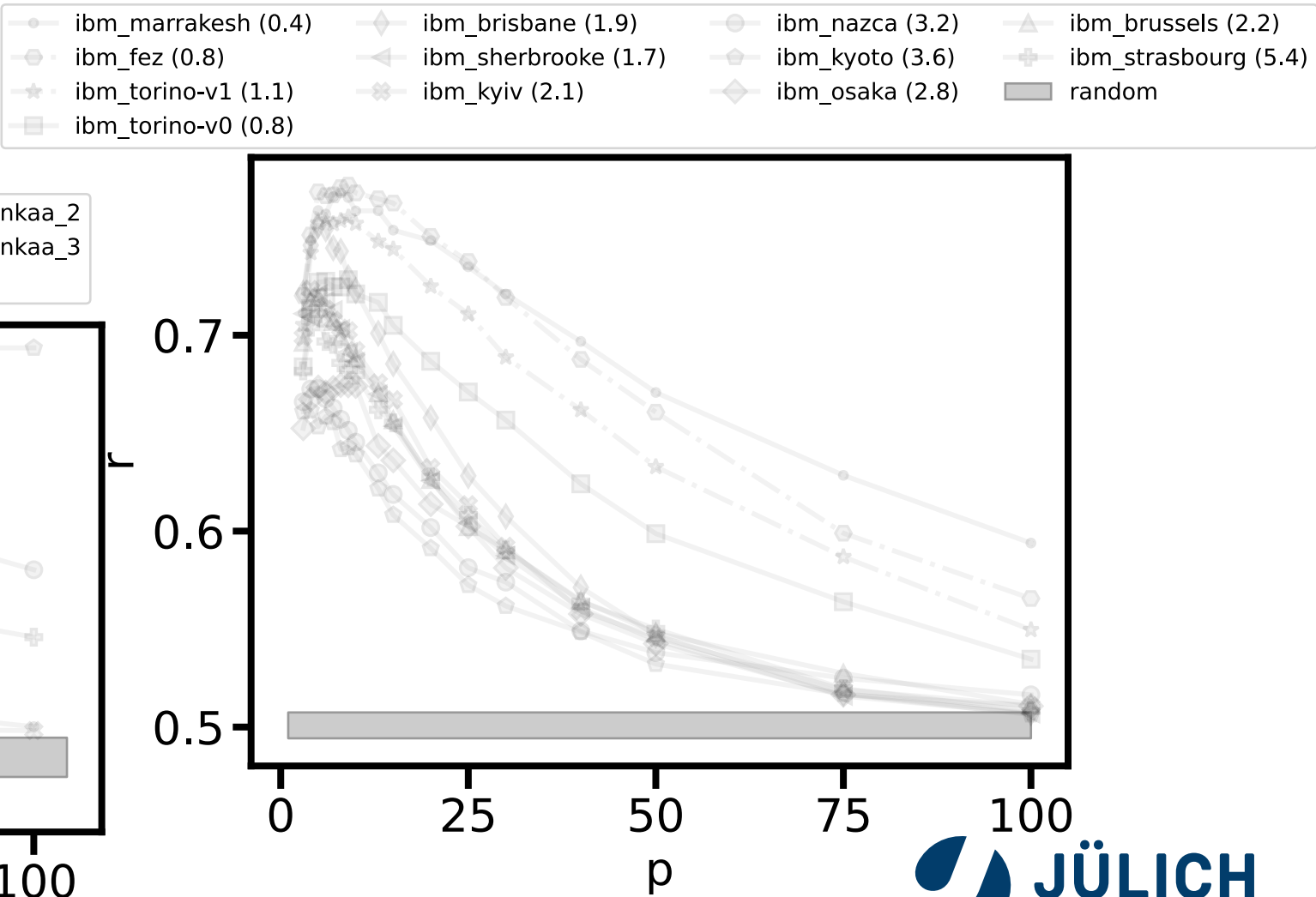
Member of the Helmholtz Association

JÜLICH Forschungszentrum

# LR-QAOA on a 1D-Chain graph



Performance on a 100-qubit 1D chain experiment. IBM QPUs (EPLG)

Legend:
- ibm_marrakesh (0.4)
- ibm_fez (0.8)
- ibm_torino-v1 (1.1)
- ibm_torino-v0 (0.8)
- ibm_brisbane (1.9)
- ibm_sherbrooke (1.7)
- ibm_kyiv (2.1)
- ibm_nazca (3.2)
- ibm_kyoto (3.6)
- ibm_osaka (2.8)
- ibm_brussels (2.2)
- ibm_strasbourg (5.4)
- random

Performance on the best 5-qubit 1D chain experiment, different QPUs

Legend:
- originq_wukong
- qasm_simulator
- iqm_garnet
- ibm_fez
- ibm_marrakesh
- ibm_brisbane
- rigetti_ankaa_2
- rigetti_ankaa_3
- random

JÜLICH
Forschungszentrum

# LR-QAOA on a 1D-Chain graph

Performance on a 100-qubit 1D chain experiment. IBM QPUs (EPLG)

Performance on the best 5-qubit 1D chain experiment, different QPUs



Member of the Helmholtz Association

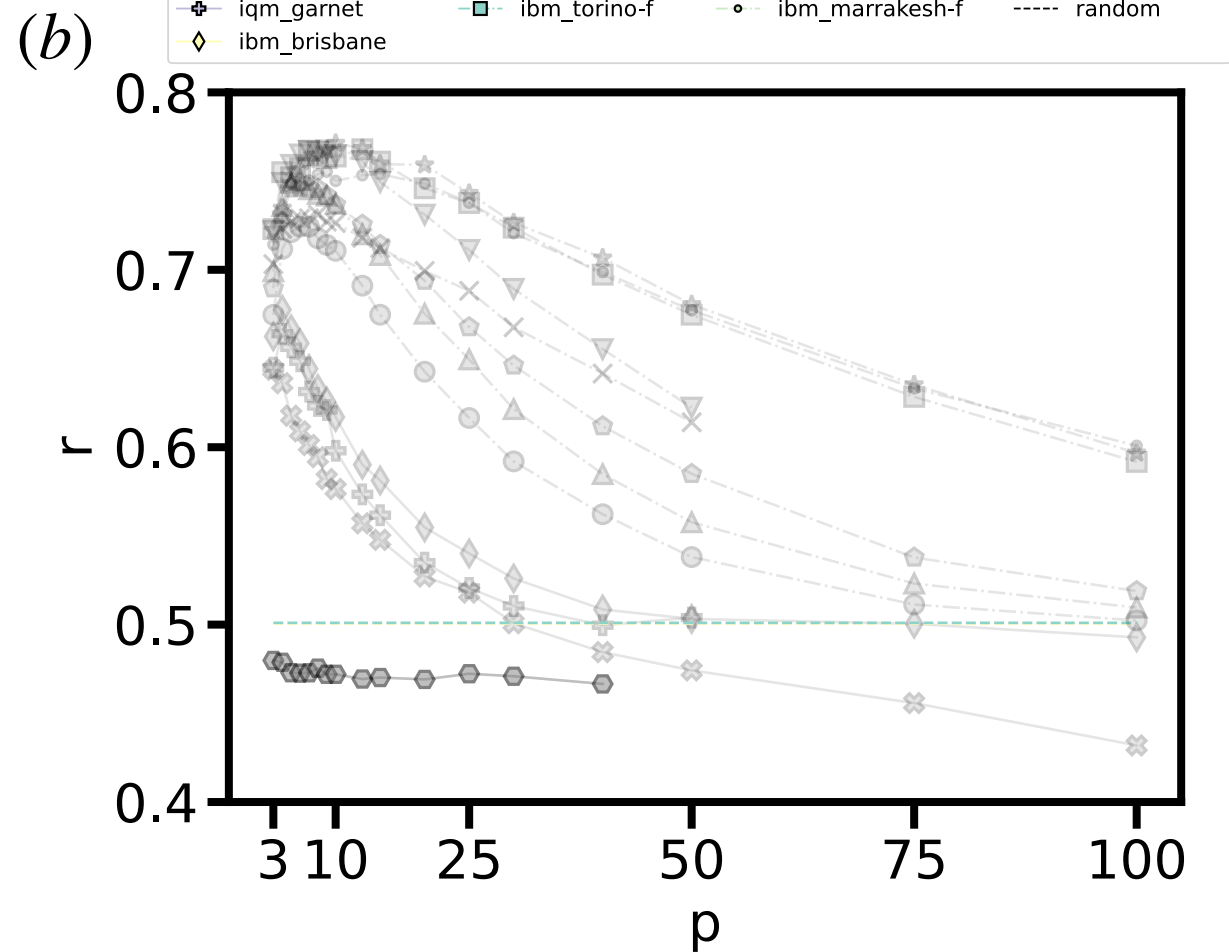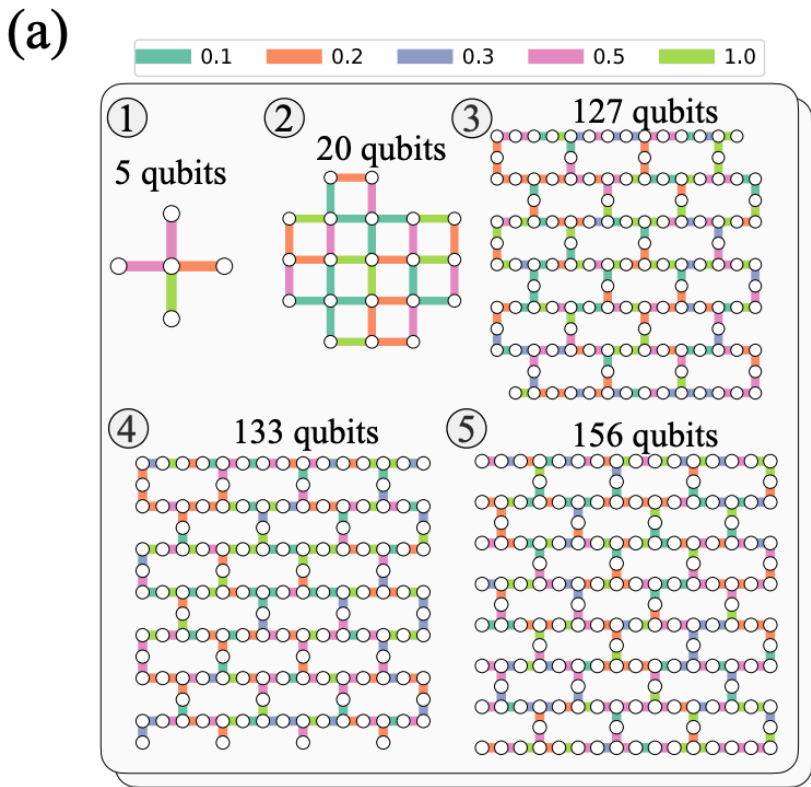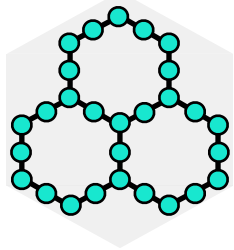JÜLICH
Forschungszentrum

# LR-QAOA on a 1D-Chain graph



Performance on the best 5-qubit 1D chain experiment, different QPUs

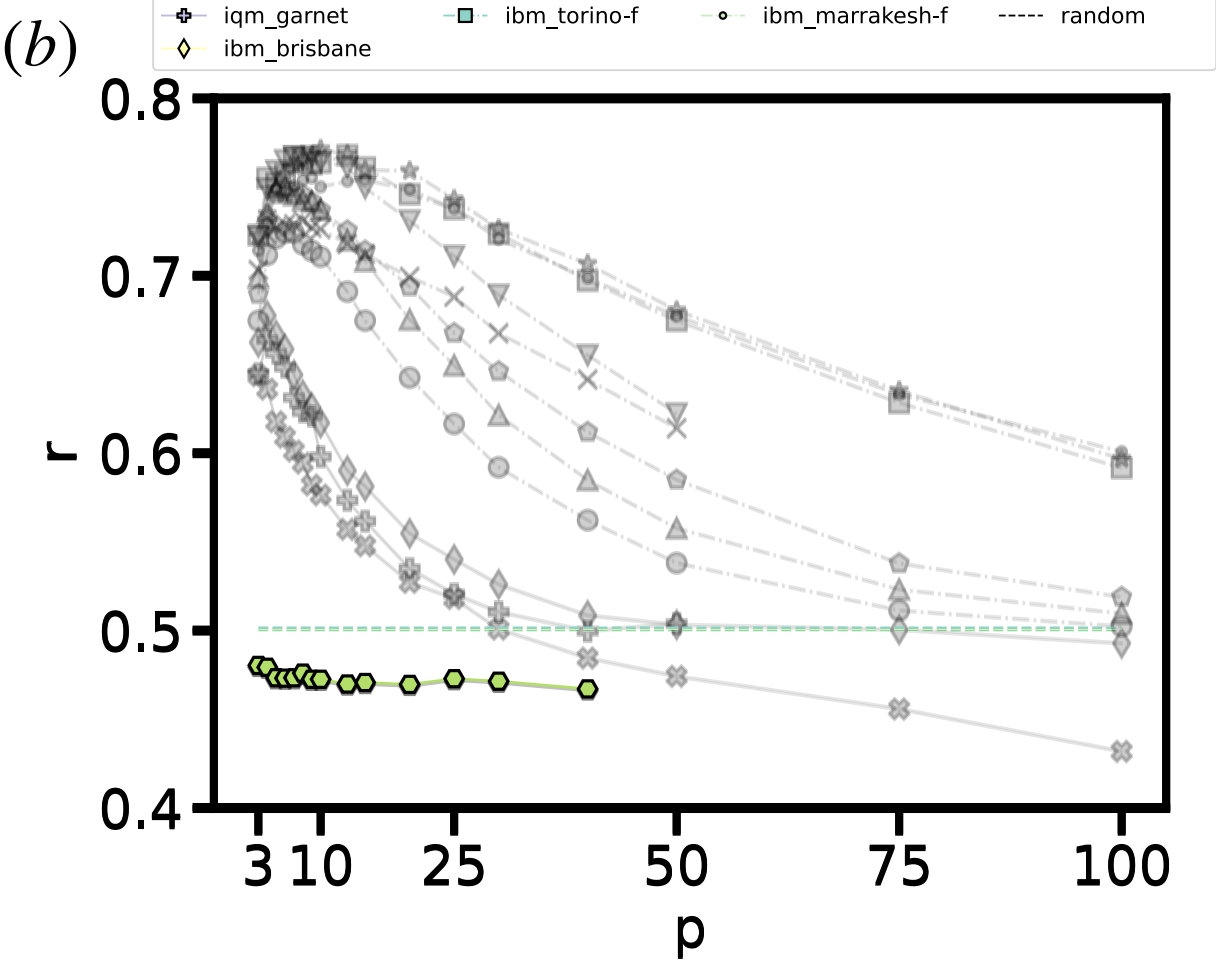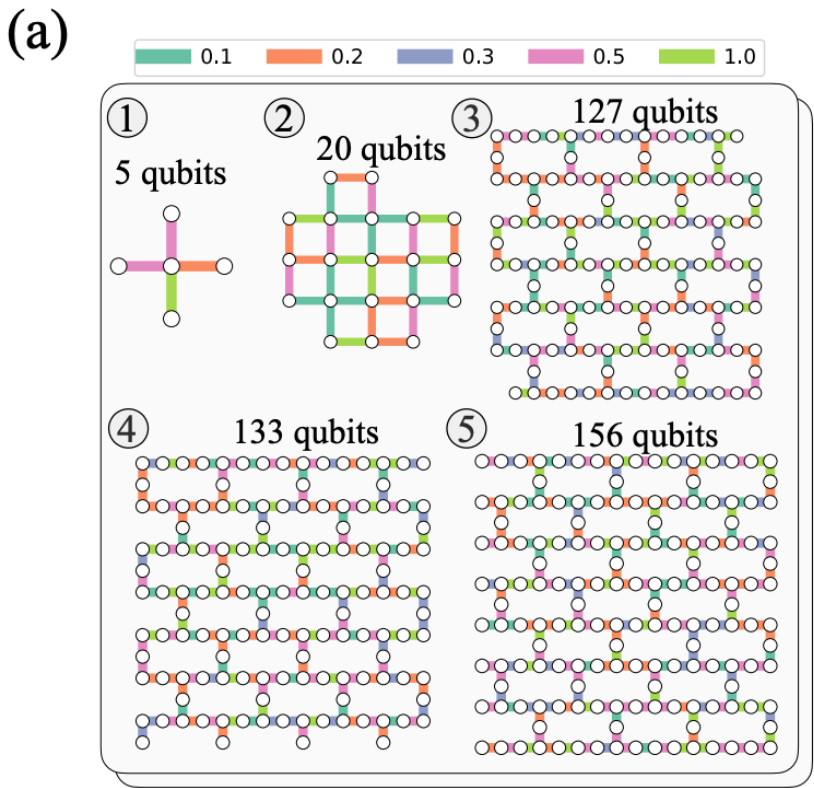Performance on a 100-qubit 1D chain experiment. IBM QPUs (EPLG)

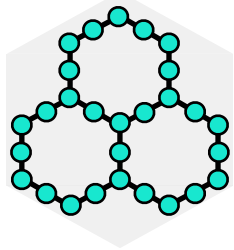Member of the Helmholtz Association

JÜLICH Forschungszentrum

# LR-QAOA on a 1D-Chain graph



Performance on a 100-qubit 1D chain experiment. IBM QPUs (EPLG)

- ibm_marrakesh (0.4)
- ibm_fez (0.8)
- ibm_torino-v1 (1.1)
- ibm_torino-v0 (0.8)
- ibm_brisbane (1.9)
- ibm_sherbrooke (1.7)
- ibm_kyiv (2.1)
- ibm_nazca (3.2)
- ibm_kyoto (3.6)
- ibm_osaka (2.8)
- ibm_brussels (2.2)
- ibm_strasbourg (5.4)
- random

Performance on the best 5-qubit 1D chain experiment, different QPUs

- originq_wukong
- qasm_simulator
- iqm_garnet
- ibm_fez
- ibm_marrakesh
- ibm_brisbane
- rigetti_ankaa_2
- rigetti_ankaa_3
- random

Member of the Helmholtz Association

JÜLICH Forschungszentrum

# LR-QAOA on a 1D-Chain graph



Performance on the best 5-qubit 1D chain experiment, different QPUs

Performance on a 100-qubit 1D chain experiment. IBM QPUs (EPLG)

# LR-QAOA on a 1D-Chain graph



Performance on a 100-qubit 1D chain experiment. IBM QPUs (EPLG)

Legend:
- ibm_marrakesh (0.4)
- ibm_fez (0.8)
- ibm_torino-v1 (1.1)
- ibm_torino-v0 (0.8)
- ibm_brisbane (1.9)
- ibm_sherbrooke (1.7)
- ibm_kyiv (2.1)
- ibm_nazca (3.2)
- ibm_kyoto (3.6)
- ibm_osaka (2.8)
- ibm_brussels (2.2)
- ibm_strasbourg (5.4)
- random

Performance on the best 5-qubit 1D chain experiment, different QPUs

Legend:
- originq_wukong
- qasm_simulator
- iqm_garnet
- ibm_fez
- ibm_marrakesh
- ibm_brisbane
- rigetti_ankaa_2
- rigetti_ankaa_3
- random

JÜLICH
Forschungszentrum

# LR-QAOA on a 1D-Chain graph



Performance on a 100-qubit 1D chain experiment. IBM QPUs (EPLG)

Performance on the best 5-qubit 1D chain experiment, different QPUs

Member of the Helmholtz Association

# LR-QAOA on a 1D-Chain graph



Performance on a 100-qubit 1D chain experiment. IBM QPUs (EPLG)

Performance on the best 5-qubit 1D chain experiment, different QPUs

# LR-QAOA on a 1D-Chain graph



Performance on a 100-qubit 1D chain experiment. IBM QPUs (EPLG)

Performance on the best 5-qubit 1D chain experiment, different QPUs

# LR-QAOA on a 1D-Chain graph



Performance on a 100-qubit 1D chain experiment. IBM QPUs (EPLG)

Performance on the best 5-qubit 1D chain experiment, different QPUs

Legend (right plot): ibm_marrakesh (0.4), ibm_fez (0.8), ibm_torino-v1 (1.1), ibm_torino-v0 (0.8), ibm_brisbane (1.9), ibm_sherbrooke (1.7), ibm_kyiv (2.1), ibm_nazca (3.2), ibm_kyoto (3.6), ibm_osaka (2.8), ibm_brussels (2.2), ibm_strasbourg (5.4), random

Legend (left plot): originq_wukong, qasm_simulator, iqm_garnet, ibm_fez, ibm_marrakesh, ibm_brisbane, rigetti_ankaa_2, rigetti_ankaa_3, random

JÜLICH Forschungszentrum

# LR-QAOA on a Native layout graph



(a) Different QPUs topologies
(b) Performance on different devices

# LR-QAOA on a Native layout graph



(a) Different QPUs topologies
(b) Performance on different devices

# LR-QAOA on a Native layout graph



(a) Different QPUs topologies
(b) Performance on different devices

# LR-QAOA on a Native layout graph



(a) Different QPUs topologies
(b) Performance on different devices

# LR-QAOA on a Native layout graph



(a) Different QPUs topologies
(b) Performance on different devices

# LR-QAOA on a Native layout graph



(a) Different QPUs topologies
(b) Performance on different devices

# LR-QAOA on a Native layout graph



(a) Different QPUs topologies
(b) Performance on different devices

# From a fixed layout to a fully connected QPUs

## SWAP networks



## Parity Twine Chain



Using a swap strategy we can convert a 1D-Chain graph into a fully connected graph. We need 3 times more 2-qubit gates to implement this protocol.

Schematic representation of the PTC encoding. (a) Circuit model to get different parities. (b) Graphic representation of the PTC in a 1D chain. Using this diagram, the equivalent building blocks for the CNOT gates are shown at the right.

JÜLICH
Forschungszentrum

# LR-QAOA on Fully Connected (FC) problems



(a) FC for a 15-qubit Weighted MaxCut problem
(b) Effective approximation ratio

$$r_{\text{eff}} = \frac{r_{\text{max}} - r_{\text{rand}}}{1 - r_{\text{rand}}}$$
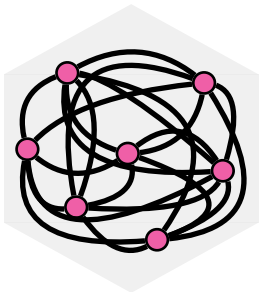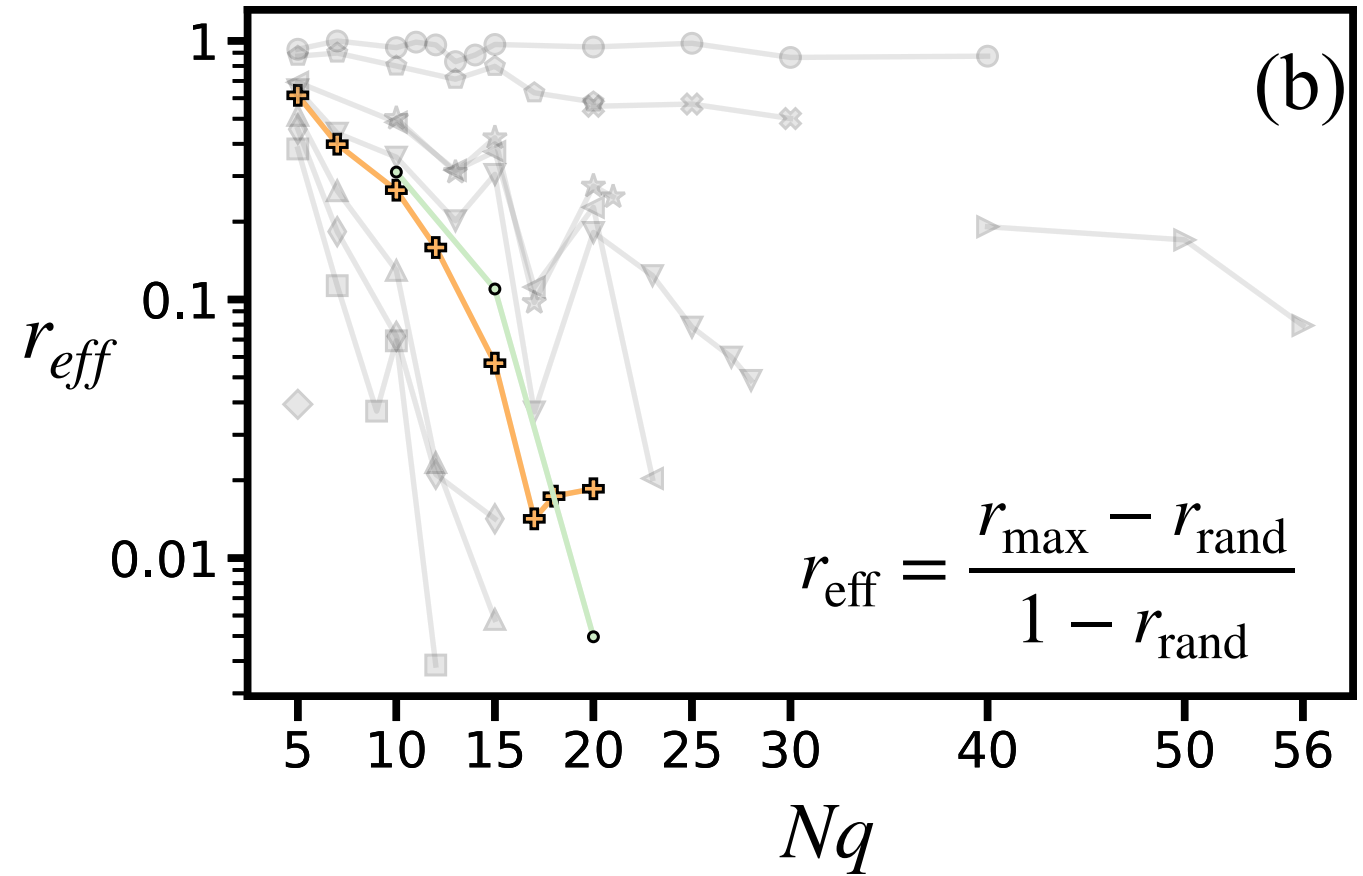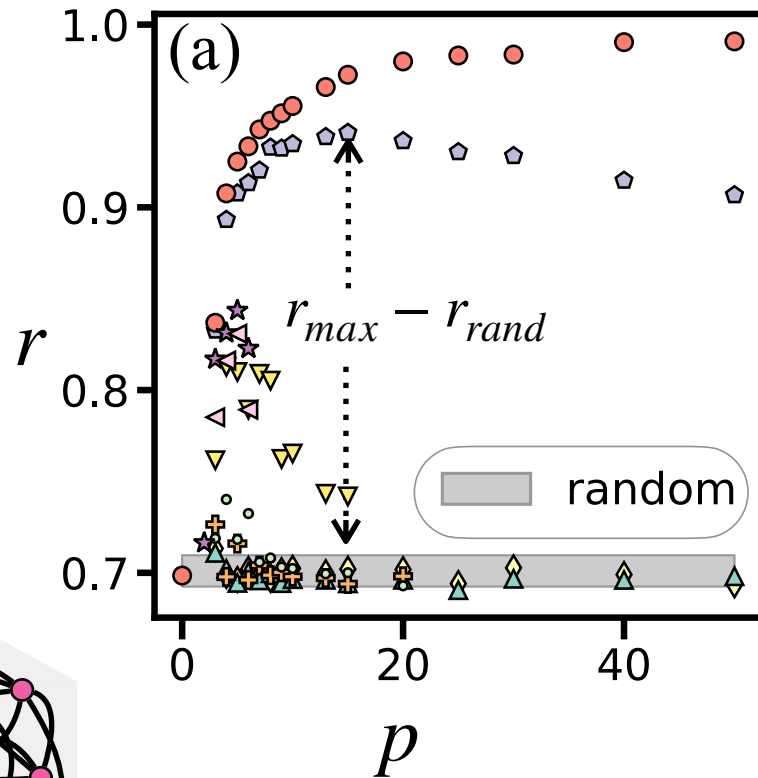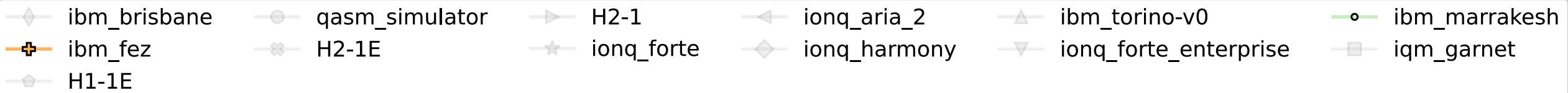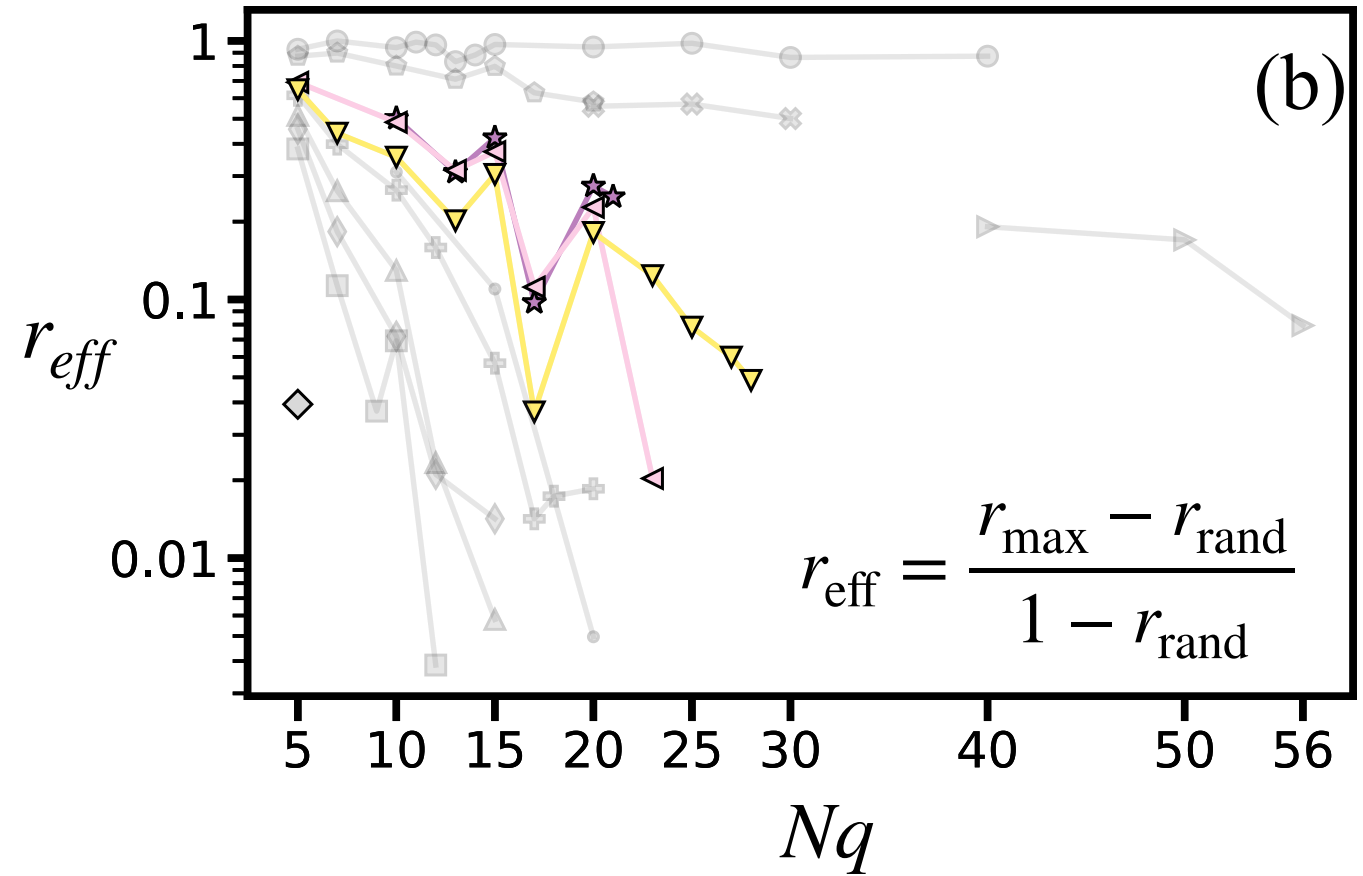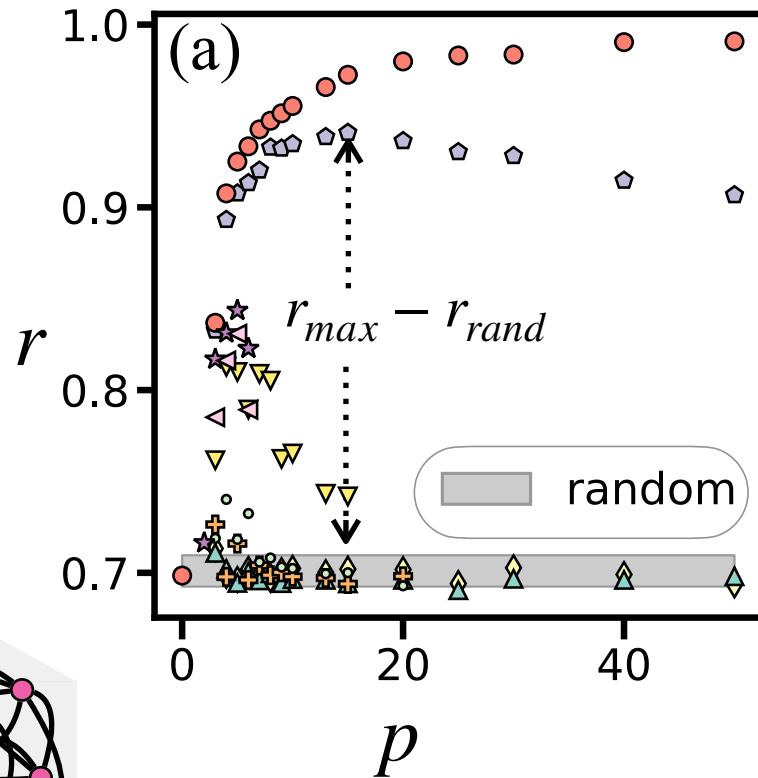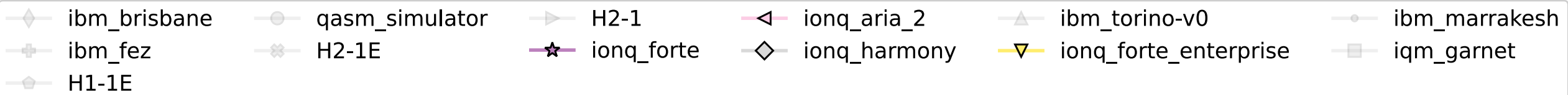
# LR-QAOA on Fully Connected (FC) problems



(a) FC for a 15-qubit Weighted MaxCut problem
(b) Effective approximation ratio

# LR-QAOA on Fully Connected (FC) problems



(a) FC for a 15-qubit Weighted MaxCut problem
(b) Effective approximation ratio
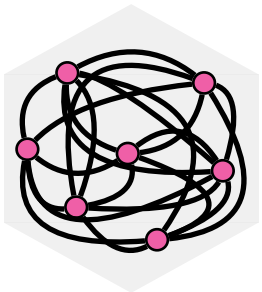
# LR-QAOA on Fully Connected (FC) problems



(a) FC for a 15-qubit Weighted MaxCut problem
(b) Effective approximation ratio

# LR-QAOA on Fully Connected (FC) problems



(a) FC for a 15-qubit Weighted MaxCut problem
(b) Effective approximation ratio

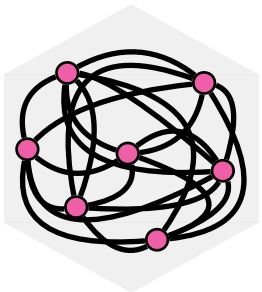# LR-QAOA on Fully Connected (FC) problems



(a) FC for a 15-qubit Weighted MaxCut problem
(b) Effective approximation ratio

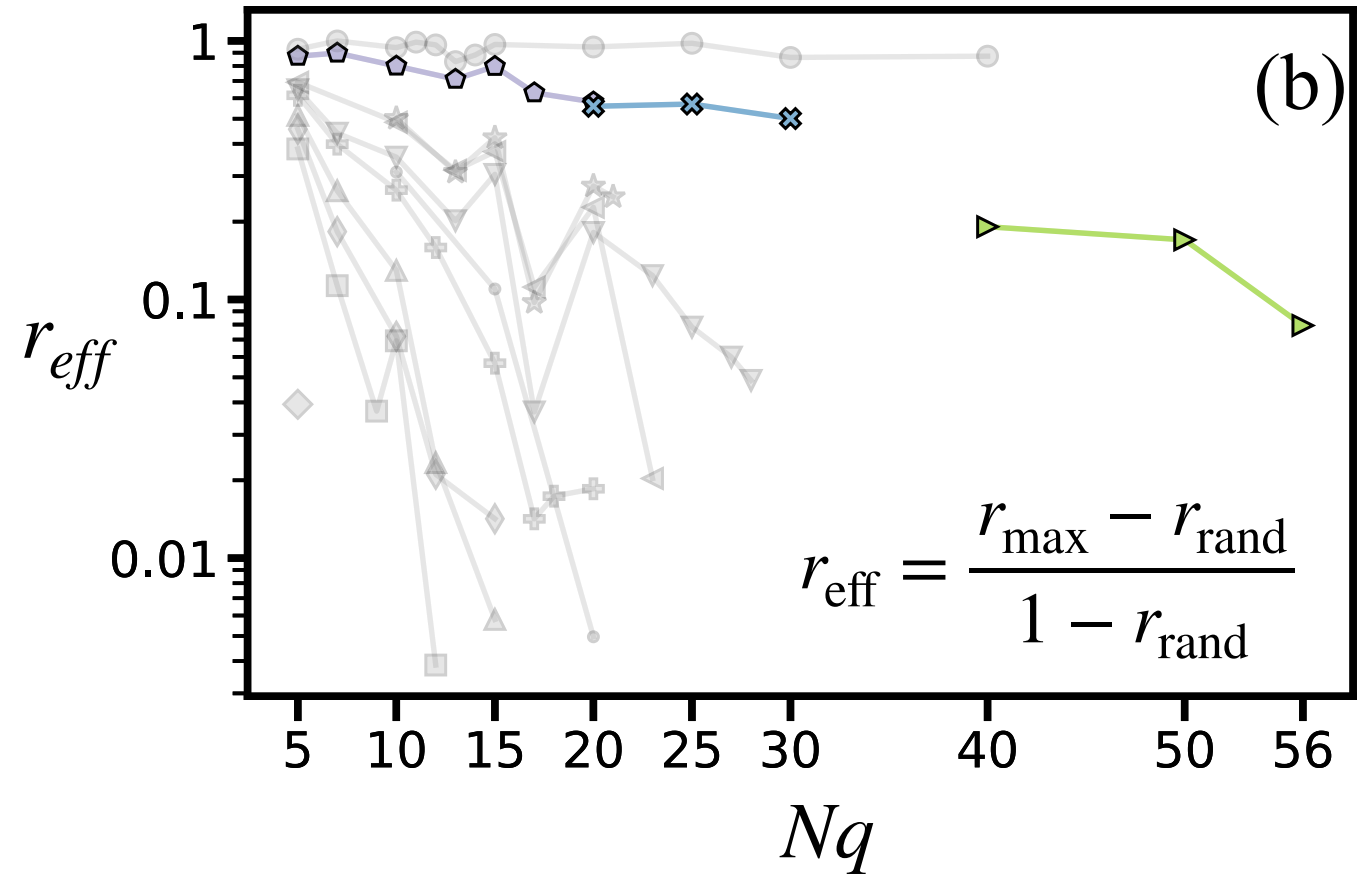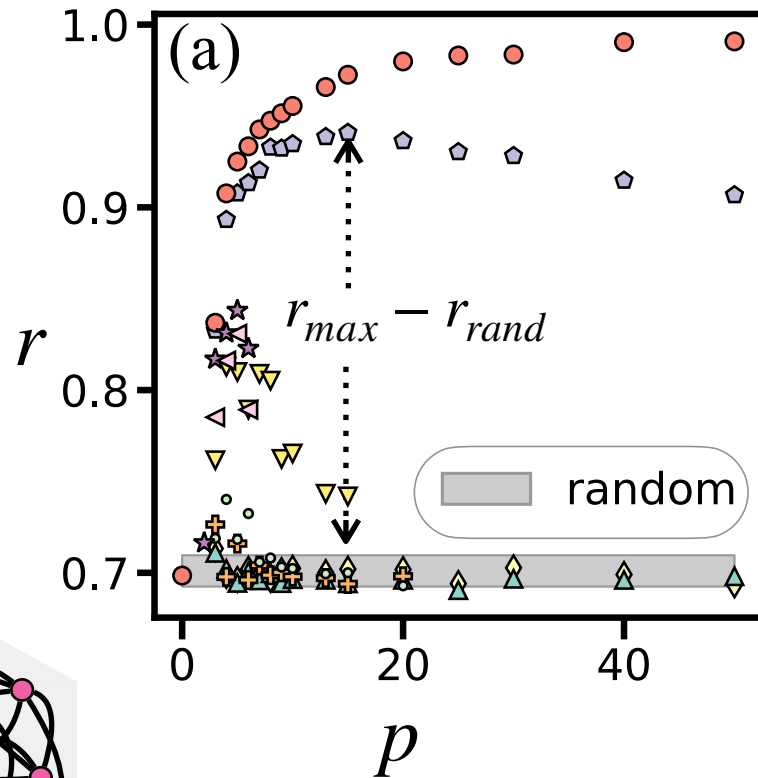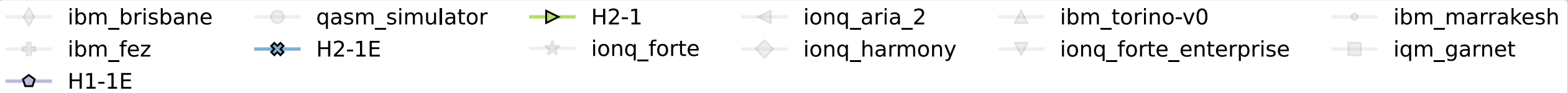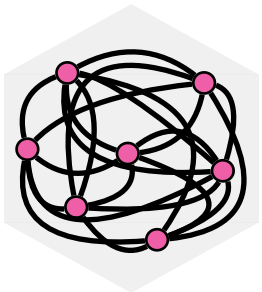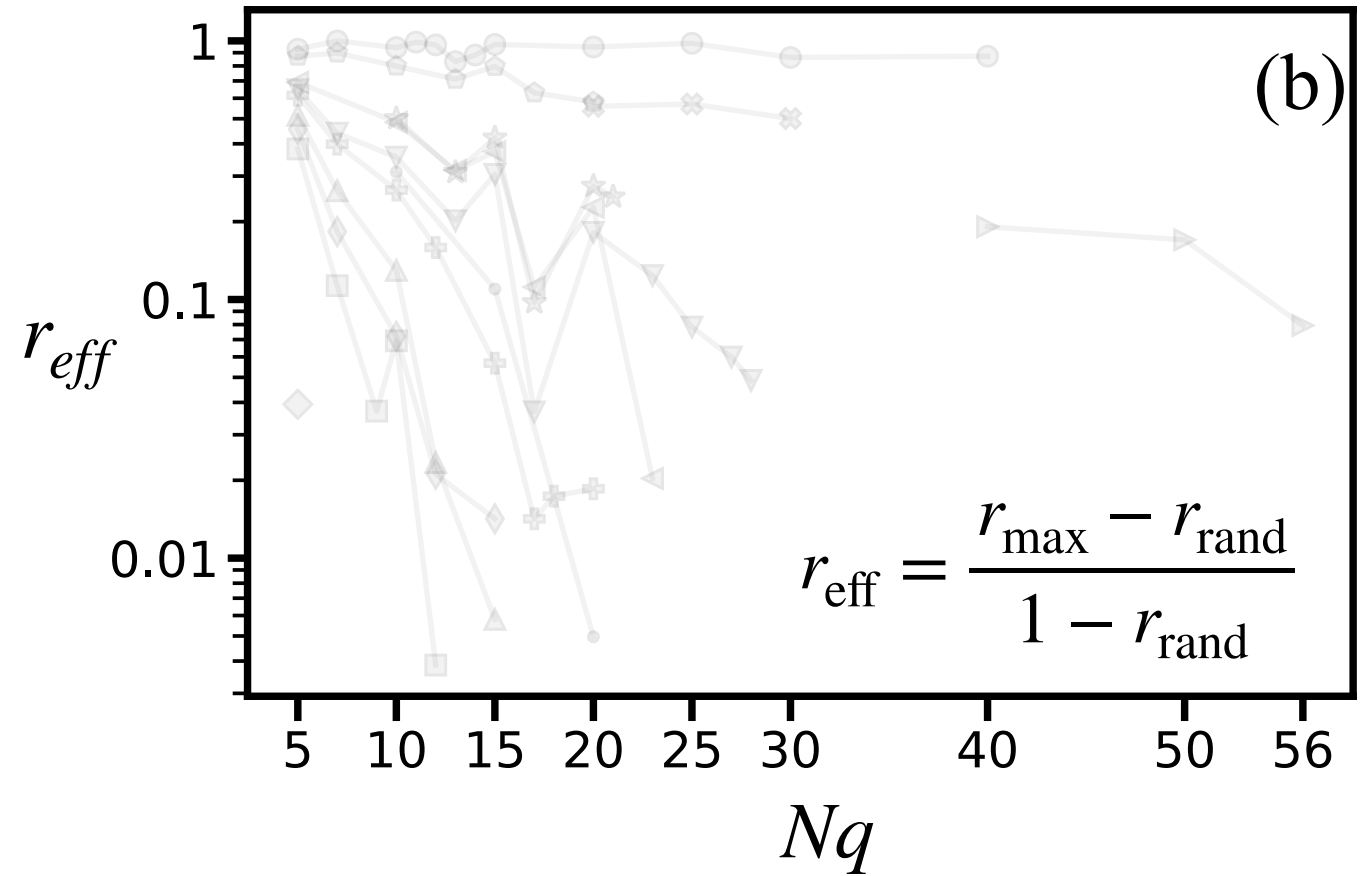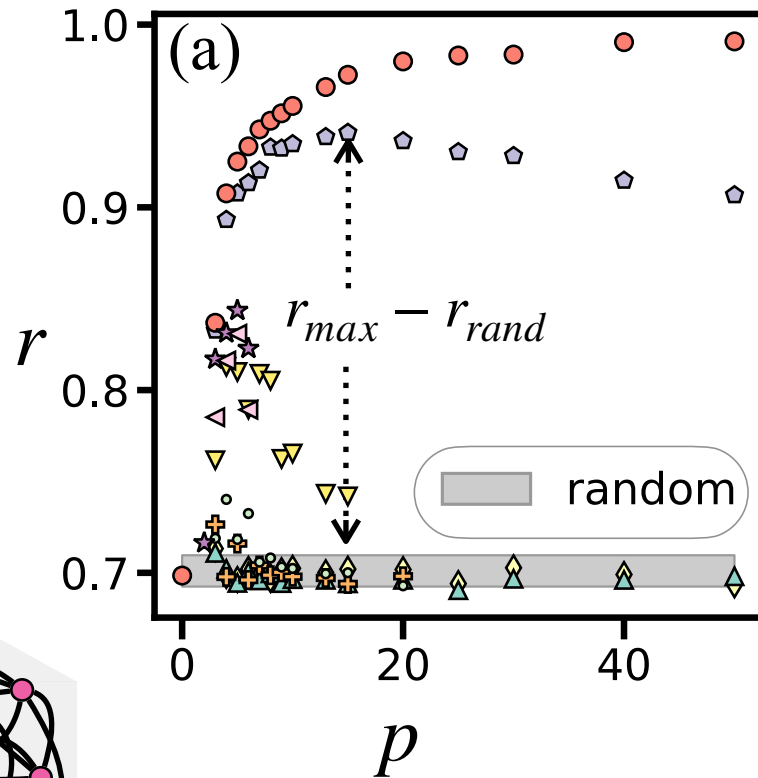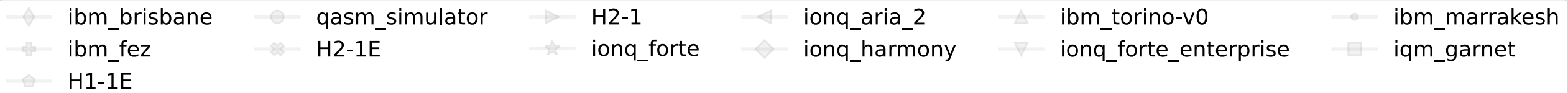# LR-QAOA on Fully Connected (FC) problems



(a) FC for a 15-qubit Weighted MaxCut problem
(b) Effective approximation ratio

# Integration LR-QAOA as an open source benchmark



**metriq-gym**

Supported By Unitary Foundation | Discord 376 online. | Contributor Covenant 2.1

`metriq-gym` is a Python framework for implementing and running standard quantum benchmarks on different quantum devices by different providers.

- *Open* – Open-source since its inception and fully developed in public.
- *Transparent* – All benchmark parameters are defined in a schema file and the benchmark code is reviewable by the community.
- *Cross-platform* – Supports running benchmarks on multiple quantum hardware providers (*integration powered by* qBraid-SDK)
- *User-friendly* – Provides a simple command-line interface for dispatching, monitoring, and polling benchmark jobs (you can go on with your life while your job waits in the queue).

JÜLICH
Forschungszentrum

# Conclusions

- We holistically benchmarked 24 QPUs from five vendors using LR-QAOA, evaluating their performance on different graph topologies and testing scalability in qubit count and circuit depth.
- IBM QPUs show significant improvements from Eagle to Heron generations, while IonQ and Quantinuum maintain performance through generations and offer better gate fidelity but suffer from slow execution times.
- Our results highlight key bottlenecks in quantum hardware, emphasizing the need for advancements in circuit depth, execution speed, and gate fidelity to support large-scale quantum algorithms.

JÜLICH
Forschungszentrum

# Thank you

JÜLICH
Forschungszentrum