

# PyQBench: a Python library for benchmarking gate-based quantum computers

Konrad Jałowiecki\*, Paulina Lewandowska, Łukasz Paweł

*Institute of Theoretical and Applied Informatics, Polish Academy of Sciences,  
Balttycka 5, 44-100 Gliwice, Poland*

---

## Abstract

We introduce PyQBench, an innovative open-source framework for benchmarking gate-based quantum computers. PyQBench can benchmark NISQ devices by verifying their capability of discriminating between two von Neumann measurements. PyQBench offers a simplified, ready-to-use, command line interface (CLI) for running benchmarks using a predefined family of measurements. For more advanced scenarios, PyQBench offers a way of employing user-defined measurements instead of predefined ones.

*Keywords:* Quantum computing, Benchmarking quantum computers, Discrimination of quantum measurements, Discrimination of von Neumann measurements, Open-source, Python programming

*PACS:* 03.67.-a, 03.67.Lx

*2000 MSC:* 81P68

---

\*Corresponding author

*Email address:* `dexter2206@gmail.com` (Konrad Jałowiecki)

## Current code version

C1	Current code version	0.1.1
C2	Permanent link to code/repository used for this code version	<a href="https://github.com/iitis/PyQBench">https://github.com/iitis/PyQBench</a>
C3	Code Ocean compute capsule	<a href="https://codeocean.com/capsule/89088992-9a27-4712-8525-d92a9b23060f/tree">https://codeocean.com/capsule/89088992-9a27-4712-8525-d92a9b23060f/tree</a>
C4	Legal Code License	Apache License 2.0
C5	Code versioning system used	git
C6	Software code languages, tools, and services used	Python, Qiskit, AWS Braket
C7	Compilation requirements, operating environments & dependencies	Python >= 3.8 numpy ~= 1.22.0 scipy ~= 1.7.0 pandas ~= 1.5.0 amazon-braket-sdk >= 1.11.1 pydantic ~= 1.9.1 qiskit ~= 0.37.2 mthree ~= 1.1.0 tqdm ~= 4.64.1 pyyaml ~= 6.0 qiskit-braket-provider ~= 0.0.3
C8	If available Link to developer documentation/manual	<a href="https://pyqbench.readthedocs.io/en/latest/">https://pyqbench.readthedocs.io/en/latest/</a>
C9	Support email for questions	dexter2206@gmail.com

Table 1: Code metadata

## 1. Motivation and significance

Noisy Intermediate-Scale Quantum (NISQ) [1] devices are storming the market, with a wide selection of devices based on different architectures and accompanying software solutions. Among hardware providers offering public access to their gate-based devices, one could mention Rigetti [2], IBM [3], Oxford Quantum Group [4], IonQ [5] or Xanadu [6]. Other vendors offer devices operating in different paradigms. ~~Notably, one could mention e.g.~~ D-Wave [7] and their quantum annealers, or QuEra devices [8] based on neutral atoms. Most vendors provide their own software stack ~~and application programming interface for accessing for interacting with~~ their devices. ~~To name a few,~~ Rigetti’s computers are available through ~~their~~ Forest SDK [9] and PyQuil

library [10] and IBM Q [3] computers can be accessed through Qiskit [11] or IBM Quantum Experience web interface [12]. ~~Some cloud services, like~~ Amazon Braket [13] ~~, offer offers~~ access to several quantum devices under a unified API. ~~On top of that, several libraries and frameworks can~~ There are also several libraries able to integrate with multiple hardware vendors. Examples ~~of such frameworks~~ include IBM Q’s Qiskit, Zapata Computing’s Orchestra [14], XACC [15] and NVIDIA’s CUDA Quantum [16].

It is well known that NISQ devices have their limitations [17]. The question is to what extent those devices can perform meaningful computations? To answer this question, one has to devise a methodology for benchmarking them. For gate-based computers, on which this paper focuses, there already exist several approaches. One could mention randomized benchmarking [18, 19, 20, 21, 22], benchmarks based on the quantum volume [23, 24, 25] or quantum tomography [26].

In this paper, we introduce a different approach to benchmarking gate-based devices with a simple operational interpretation. In our method, we test how ~~well~~ good the given device is at guessing which of the two known von Neumann measurements were performed during the experiment. We implemented our approach in an open-source Python library ~~called~~ PyQBench. The library supports any device available through the Qiskit library, ~~and thus can be used with providers such as~~ including all IBM Q or Amazon Braket. Along with the library, the PyQBench package contains a command line tool for running the most common benchmarking scenarios.

## ~~2. Existing benchmarking methodologies and software~~

Unsurprisingly, PyQBench is not the only software package for benchmarking gate-based devices. ~~While we believe that our approach has significant benefits over other benchmarking techniques, for completeness, in this section we discuss some of the currently available similar software.~~

~~Probably the simplest benchmarking method one could devise is simply running known algorithms and comparing outputs with the expected ones. Analyzing the frequency of the correct outputs, or the deviation between actual and expected outputs distribution provides then a metric of the performance of a given device. Libraries such as Munich Quantum Toolkit (MQT) [27, 28] or SupermarQ [29, 30] contain benchmarks leveraging multiple algorithms, such as Shor’s algorithm or Grover’s algorithm. Despite being intuitive and easily interpretable, such benchmarks may have some problems. Most importantly, they assess the usefulness of a quantum device only for a very particular algorithm, and it might be hard to extrapolate their results to other algorithms and applications. For instance, the inability of a device to~~

consistently find factorizations using Shor’s algorithms does not tell anything about its usefulness in Variational Quantum Algorithm’s.

Another possible approach to benchmarking quantum computers is randomized benchmarking. In this approach, one samples circuits to be run from some predefined set of gates (e.g. from the Clifford group) and tests how much the output distribution obtained from the device running these circuits differs from the ideal one. It is also common to concatenate randomly chosen circuits with their inverses (which should yield the identity circuit) and run those concatenated circuits on the device. Libraries implementing this approach include Qiskit [31] or PyQuil [32]. Another, equally popular, benchmarking method is quantum tomography [26]. Implementation of quantum tomography for benchmarking NISQ devices can be found in [33].

In [34], the authors evaluated several IBM-Q machines using seven benchmarks taking into account the errors and execution time.

QASMBench [35] is one of the first benchmark suites aiming at evaluating NISQ devices using quantum applications from a broad range of domains, mainly using an approach based on fidelity estimation. The benchmark presented in the paper compares the fidelity of execution among the IBM-Q machines, the IonQ QPU and the Rigetti Aspen M-1 system.

Another quantity used for benchmarking NISQ devices is quantum volume. The quantum volume characterizes the capacity of a device for solving computational problems. It takes into account multiple factors like the number of qubits, connectivity and measurement errors. The Qiskit library allows one to measure the quantum volume of a device by using its `qiskit.ignis.verification.quantum-volume`. Other implementations of Quantum Volume can be found as well, see e.g. [36].

We should also mention cross-entropy benchmarking [37], which was utilized in validation of the Sycamore-53 QPU supremacy experiments [38]. In this approach, the quality of an algorithm implemented on the QPU is measured by calculating the cross-entropy of bit-strings actually sampled from the QPU, compared to ideal bitstrings.

Finally, it is worth pointing out there is an ongoing effort towards standardization of benchmarking of quantum computers. In [39], the authors present plans for designing a benchmarking suite based on measuring a set of standardized key performance indicators (KPIs).

## 2. Preliminaries and discrimination scheme approach

In this section, we describe how the benchmarking process in PyQBench works. We start by discussing necessary mathematical preliminaries. Then, we present the general form of the discrimination scheme used in PyQBench

90 and ~~practical considerations on how to implement it taking into account~~  
 91 ~~the limitations of the current~~ discuss it in the context of current, limited  
 92 NISQ devices. We ~~encourage the refer~~ readers interested in ~~a more in-depth~~  
 93 ~~discussion of more details concerning~~ the mathematical foundations behind  
 94 ~~our discrimination scheme to read~~ PyQBench to Section 2 in the ~~supplemental~~  
 95 ~~materials~~ supplementary material.

### 96 2.1. Von Neumann Measurements

97 A von Neumann measurement  $\mathcal{P}$  is a collection of rank-one projectors  
 98  $\{|u_0\rangle\langle u_0|, \dots, |u_{d-1}\rangle\langle u_{d-1}|\}$  ~~called effects~~, that sum up to the identity op-  
 99 erator, i.e.  $\sum_{i=0}^{d-1} |u_i\rangle\langle u_i| = \mathbb{1}$ . If  $U$  is a unitary matrix of size  $d$ , one can  
 100 construct a von Neumann measurement  $\mathcal{P}_U$  by taking projectors onto its  
 101 columns. In this case, we say that  $\mathcal{P}_U$  is described by the matrix  $U$ .

102 Typically, NISQ devices can only perform measurements in computational  
 103  $Z$ -basis, i.e.  $U = \mathbb{1}$ . To implement an arbitrary von Neumann measurement  
 104  $\mathcal{P}_U$ , one has to first apply  $U^\dagger$  to the measured system and then follow with  
 105  $Z$ -basis measurement. This process, depicted in Fig. 1, can be viewed as  
 106 performing a change of basis in which measurement is performed prior to  
 107 measurement in the computational basis.

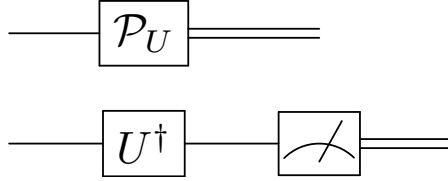


Figure 1: Implementation of a von Neumann measurement using measurement in computational basis. The upper circuit shows a symbolic representation of a von Neumann measurement  $\mathcal{P}_U$ . The bottom, equivalent circuit depicts its decomposition into a change of basis followed by measurement in the  $Z$  basis.

### 108 2.2. Discrimination scheme

109 Benchmarks in PyQBench work by experimentally determining the prob-  
 110 ability of correct discrimination between two von Neumann measurements  
 111 by the device under test and comparing the result with the ideal, theoretical  
 112 predictions.

113 Without loss of generality<sup>1</sup>, we consider discrimination task between sin-  
 114 gle qubit measurements  $\mathcal{P}_1$ , performed in the ~~computational~~  $Z$ -basis, and an

<sup>1</sup>Explaining why we can consider only discrimination scheme between  $\mathcal{P}_1$  and  $\mathcal{P}_U$  is beyond the scope of this paper. See [40] for a in depth explanation.

115 alternative measurement  $\mathcal{P}_U$  performed in the basis  $U$ . The ~~discrimination~~  
 116 scheme presented in Fig. 2 requires an auxiliary qubit. First, the joint system  
 117 is prepared in some state  $|\psi_0\rangle$ . Then, one of the measurements, either  $\mathcal{P}_U$   
 118 or  $\mathcal{P}_1$ , is performed on the first part of the system. Based on its outcome  
 119  $i$ , we choose another binary measurement  $\mathcal{P}_{V_i}$  and perform it on the second  
 120 qubit, obtaining the outcome  $j$ . Finally, if  $j = 0$ , we say that the performed  
 121 measurement is  $\mathcal{P}_U$ , otherwise we say that it was  $\mathcal{P}_1$ .

122 Note, however, that the discrimination scheme described above can work  
 123 regardless of the dimensionality of the unitary  $U$ . The main difference is  
 124 the dimension of the auxiliary system, which in general can be larger than  
 125 two. This dimension depends on the Schmidt rank of the optimal input  
 126 state. However, for most discrimination schemes, the Schmidt rank equals at  
 127 most two, and ~~hence~~ the auxiliary system is also a qubit, see [40] for details.  
 128 Note that the final measurement  $\mathcal{P}_{V_i}$  is always binary, independently of the  
 129 dimension of the auxiliary system.

130 Naturally, we need to repeat the same procedure multiple times for both  
 131 measurements to obtain a reliable estimate of the underlying probability  
 132 distribution. In PyQBench, we assume that the experiment is repeated the  
 133 same number of times for both  $\mathcal{P}_U$  and  $\mathcal{P}_1$ .

134 In principle, our discrimination scheme could be used with any choice of  
 135  $|\psi_0\rangle$  and final measurements  $\mathcal{P}_{V_i}$ . However, we argue that it is best to choose  
 136 those components so that they maximize the probability of correct discrim-  
 137 ination. To see that, suppose that some choice of  $|\psi_0\rangle, \mathcal{P}_{V_0}, \mathcal{P}_{V_1}$  allows for  
 138 correctly discriminating between two measurements with probability equal  
 139 to one, i.e. on a perfect quantum computer you will always make a cor-  
 140 rect guess. Then, on real hardware, we might obtain any empirical value in  
 141 range  $[\frac{1}{2}, 1]$ . On the other hand, if we choose the components of our scheme  
 142 such that the successful discrimination probability is  $\frac{3}{5}$ , the possible range  
 143 of empirically obtainable probabilities is only  $[\frac{1}{2}, \frac{3}{5}]$ . Hence, in the second  
 144 case, the discrepancy between theoretical and empirical results will be less  
 145 pronounced.

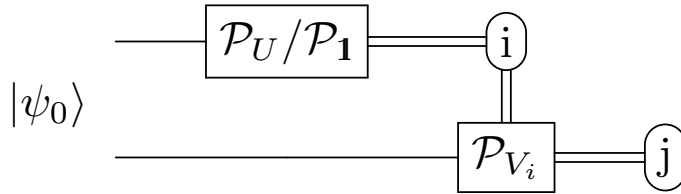


Figure 2: Theoretical scheme of discrimination between von Neumann measurements  $\mathcal{P}_U$  and  $\mathcal{P}_1$ .

146 2.2.1. Implementation of discrimination scheme on actual NISQ devices

147 Current NISQ devices are unable to perform conditional measurements,  
 148 which is the biggest obstacle to implementing our scheme on real hardware.  
 149 ~~However, we~~ We circumvent this problem by slightly adjusting our scheme  
 150 ~~so that it only uses~~ to only use components available on current devices. ~~For~~  
 151 ~~this purpose, we use~~ There are two possible options: using a postselection or  
 152 a direct sum  $V_0^\dagger \oplus V_1^\dagger$ .

153 **Scheme 1.** (Postselection)

154 The first idea uses a postselection scheme. In the original scheme, we  
 155 measure the first qubit and only then determine which measurement should  
 156 be performed on the second one. ~~Instead of doing this choice,~~ we can run two  
 157 circuits, one with  $\mathcal{P}_{V_0}$  and one with  $\mathcal{P}_{V_1}$  and measure both qubits. We can  
 158 then discard the results of the circuit for which label  $i$  does not match mea-  
 159 surement label  $k$ . ~~Hence, the~~ The circuit for postselection looks as depicted  
 160 in Fig. 3.

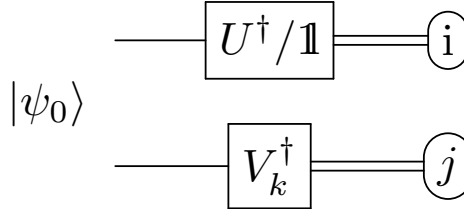


Figure 3: A schematic representation of the setup for distinguishing measurements  $\mathcal{P}_U$  and  $\mathcal{P}_1$  using postselection approach. In postselection scheme, one runs such circuits for both  $k = 0, 1$  and discards results for cases when there is a mismatch between  $k$  and  $i$ .

161 To perform the benchmark, one needs to run multiple copies of the post-  
 162 selection circuit, with both  $\mathcal{P}_U$  and  $\mathcal{P}_1$ . Each circuit has to be run in both  
 163 variants, one with final measurement  $\mathcal{P}_{V_0}$  and the second with the final mea-  
 164 surement  $\mathcal{P}_{V_1}$ . The experiments can thus be grouped into classes identified by  
 165 tuples of the form  $(\mathcal{Q}, k, i, j)$ , where  $\mathcal{Q} \in \{\mathcal{P}_U, \mathcal{P}_1\}$  denotes the chosen mea-  
 166 surement,  $k \in \{0, 1\}$  designates the final measurement used, and  $i \in \{0, 1\}$   
 167 and  $j \in \{0, 1\}$  being the labels of outcomes as presented in Fig. 3. We  
 168 then discard all the experiments for which  $i \neq k$ . The total number of valid  
 169 experiments is thus:

$$N_{\text{total}} = \#\{(\mathcal{Q}, k, i, j) : k = i\}. \quad (1)$$

170 Finally, we count the valid experiments resulting in successful discrimi-  
 171 nation. If we have chosen  $\mathcal{P}_U$ , then we guess correctly iff  $j = 0$ . Similarly,

for  $\mathcal{P}_1$ , we guess correctly iff  $j = 1$ . If we define

$$N_{\mathcal{P}_U} = \#\{(\mathcal{Q}, k, i, j) : \mathcal{Q} = \mathcal{P}_U, k = i, j = 0\}, \quad (2)$$

$$N_{\mathcal{P}_1} = \#\{(\mathcal{Q}, k, i, j) : \mathcal{Q} = \mathcal{P}_1, k = i, j = 1\}, \quad (3)$$

then the empirical success probability can be computed as

$$p_{\text{succ}}(\mathcal{P}_U, \mathcal{P}_1) = \frac{N_{\mathcal{P}_U} + N_{\mathcal{P}_1}}{N_{\text{total}}}. \quad (4)$$

The  $p_{\text{succ}}$  is the quantity reported to the user as the result of the benchmark.

**Scheme 2.** (Direct sum)

The second idea uses the direct sum  $V_0^\dagger \oplus V_1^\dagger$  implementation. Here, instead of performing a conditional measurement  $\mathcal{P}_{V_k}$ , where  $k \in \{0, 1\}$ , we run circuits presented in Fig. 4.

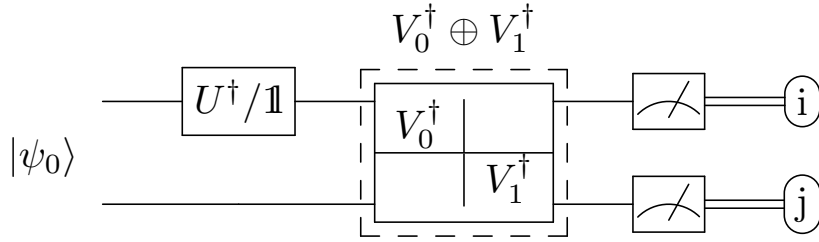


Figure 4: A schematic representation of the setup for distinguishing measurements  $\mathcal{P}_U$  and  $\mathcal{P}_1$  using the  $V_0^\dagger \oplus V_1^\dagger$  direct sum.

One can see why such a circuit is equivalent to the original discrimination scheme. If we rewrite the block-diagonal matrix  $V_0^\dagger \oplus V_1^\dagger$  as follows:

$$V_0^\dagger \oplus V_1^\dagger = |0\rangle\langle 0| \otimes V_0^\dagger + |1\rangle\langle 1| \otimes V_1^\dagger, \quad (5)$$

we can see that the direct sum in Eq. (5) commutes with the measurement on the first qubit. Thanks to this, we can switch the order of operations to obtain the circuit from Fig. 5. Now, depending on the outcome  $i$ , one of the summands in Eq. (5) vanishes, and we end up performing ~~exactly~~ the same operations as in the original scheme.

In this scheme, the experiment can be characterized by a pair  $(\mathcal{Q}, i, j)$ , where  $\mathcal{Q} = \{\mathcal{P}_U, \mathcal{P}_1\}$  and  $i, j \in \{0, 1\}$  are the output labels. The number of successful trials for  $U$  and  $\mathbb{1}$ , respectively, can be written as

$$N_{\mathcal{P}_U} = \#\{(\mathcal{Q}, i, j) : \mathcal{Q} = \mathcal{P}_U, j = 0\}, \quad (6)$$

$$N_{\mathcal{P}_1} = \#\{(\mathcal{Q}, i, j) : \mathcal{Q} = \mathcal{P}_1, j = 1\}. \quad (7)$$



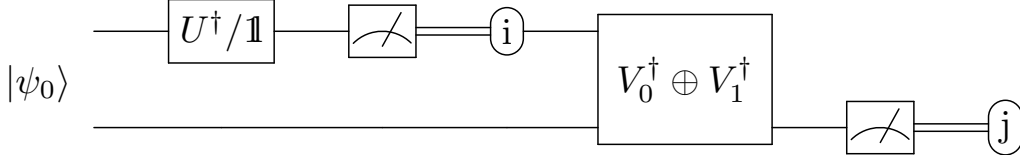


Figure 5: Rewritten representation of the setup for distinguishing measurements  $\mathcal{P}_U$  and  $\mathcal{P}_1$  using the  $V_0^\dagger \oplus V_1^\dagger$  direct sum.

Then, the probability of correct discrimination between  $\mathcal{P}_U$  and  $\mathcal{P}_1$  is given by

$$p_{\text{succ}} = \frac{N_{\mathcal{P}_U} + N_{\mathcal{P}_1}}{N_{\text{total}}}, \quad (8)$$

where  $N_{\text{total}}$  is the number of trials.

Compared to ~~these~~ other approaches, our approach allows for a very simple operational interpretation of the scheme and its results. This is especially useful for newcomers to the field, who may be put off by more complicated approaches. Another benefit is, especially for advanced users, the ability to control resources utilized during the benchmarking. We can consider here resources such as entanglement or coherence. Finally, the figure of merit we wish to calculate and to which we compare the results is fairly simple to obtain. The main downside of our approach is the exponential number of circuits we need to consider.

### 3. Software description

This section is divided into two parts. In Section 3.1 we describe functionalities of PyQBench package. Next, in Section 3.2, we give a general overview of the software architecture.

#### 3.1. Software Functionalities

The PyQBench can be used in two modes: as a Python library and as a CLI script. When used as a library, PyQBench allows the customization of discrimination scheme. The user provides a unitary matrix  $U$  defining the measurement to be discriminated, the discriminator  $|\psi_0\rangle$ , and unitaries  $V_0$  and  $V_1$  describing the final measurement. The PyQBench library provides then the following functionalities.

1. Assembling circuits for both postselection and direct sum-based discrimination schemes.
2. Executing the whole benchmarking scenario on specified backend (either real hardware or software simulator).

216 3. Interpreting the obtained outputs in terms of discrimination probab-  
217 ities.

218 Note that the execution of circuits by PyQBench is optional. ~~Instead, the~~  
219 ~~user might want to opt in~~ The user might opt-in for fine-grained control over  
220 the execution of the circuits. For instance, suppose the user wants to simulate  
221 the discrimination experiment on a noisy simulator. In such a case, they can  
222 define the necessary components and assemble the circuits using PyQBench.  
223 The circuits can then be altered, e.g. to add noise to particular gates, and  
224 then run using any Qiskit backend by the user. Finally, PyQBench can be  
225 used to interpret the measurements to obtain discrimination probability.

226 The PyQBench library also contains a readily available implementation  
227 of all necessary components needed to run discrimination experiments for  
228 parametrized Fourier family of measurements (see Section 4 in supplemen-  
229 tal material). However, if one only wishes to use this particular family of  
230 measurements in their benchmarks, then using PyQBench as a command  
231 line tool ~~might be more straightforward.~~ is a more straightforward way.  
232 The PyQBench's ~~command-line interface allows running the benchmarking~~  
233 ~~process without CLI~~ does not require writing Python code. ~~The configuration~~  
234 ~~of CLI is done~~ and is driven by YAML [41] files describing the benchmark  
235 to be performed and the description of the backend on which the benchmark  
236 should be run. The same benchmark can be used with different backends  
237 and vice versa.

### 238 3.2. Software Architecture

#### 239 3.2.1. Overview of the software structure

240 ~~As already described,~~ PyQBench can be used both as a library and a CLI.  
241 Both functionalities are implemented as a part of `qbench` Python package.  
242 The exposed CLI tool is also named `qbench`. For brevity, we ~~do not discuss~~  
243 ~~the exact structure of the package here, and instead limit ourselves~~ limit  
244 ourselves here to summarizing the architecture on the diagram in Fig. 6. For  
245 further details, we refer an interested reader to the source code available at  
246 GitHub [42] or at the reference manual [43].

247 PyQBench can be installed from the official Python Package Index (PyPI)  
248 by running `pip install pyqbench`. In a properly configured Python envi-  
249 ronment, the installation process should also make the `qbench` command  
250 available to the user without a need for further configuration.

#### 251 3.2.2. Integration with hardware providers and software simulators

252 PyQBench is built around the Qiskit [11] ecosystem. Hence, both the  
253 CLI tool and the `qbench` library can use any Qiskit-compatible backend.

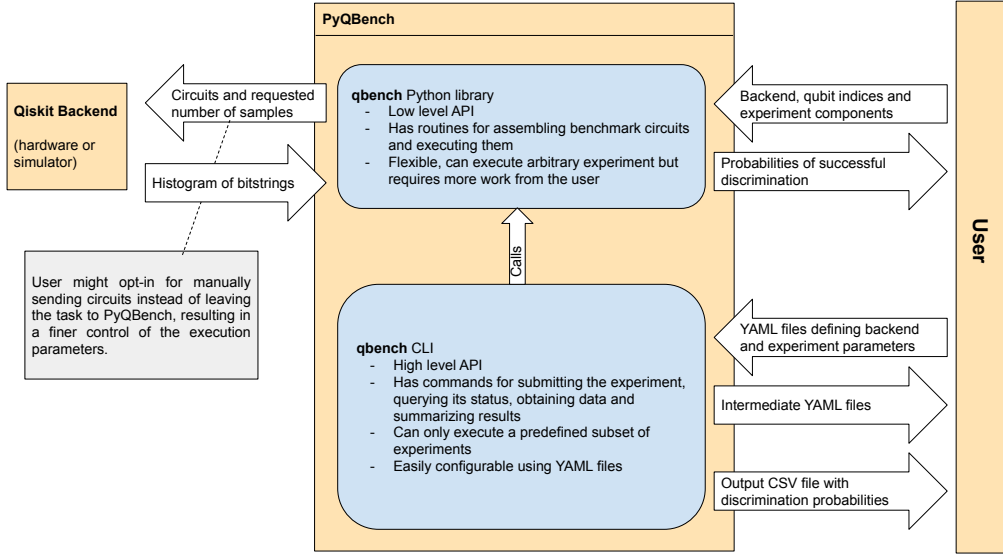


Figure 6: Overview of PyQBench's architecture.

254 This includes, IBM Q backends (available by default in Qiskit) and Amazon  
 255 Braket devices and simulators (available through `qiskit-braket-provider`  
 256 package [44, 45]).

257 When using PyQBench as a library, instances of Qiskit backends can be  
 258 passed to functions that expect them as parameters. ~~However, in~~ In CLI  
 259 mode, the user has to provide a YAML file describing the backend. An  
 260 example of such a file can be found in Section 4, and the detailed description  
 261 of the expected format can be found ~~at~~ in PyQBench's documentation.

### 262 3.2.3. Command Line Interface

263 The Command Line Interface (CLI) of PyQBench has a nested structure.  
 264 The general form of the CLI invocation is shown in listing 1.

Listing 1: Invocation of `qbench` script

```
265 qbench <benchmark-type> <command> <parameters>
```

268 Currently, PyQBench's CLI supports only one type of benchmark (discrimi-  
 269 nation of parametrized Fourier family of measurements), but we decided on  
 270 hierarchically structuring the CL to allow for future extensions. Thus, the  
 271 only accepted value of `<benchmark-type>` is `disc-fourier`. The `qbench`  
 272 `disc-fourier` command has four subcommands:

- 273 • **benchmark**: run benchmarks. This creates either a result YAML file

274 containing the measurements or an intermediate YAML file for asyn-  
275 chronous experiments.

276 • **status**: query status of experiments submitted for given benchmark.  
277 This command is only valid for asynchronous experiments.

278 • **resolve**: query the results of asynchronously submitted experiments  
279 and write the result YAML file. The output of this command is almost  
280 identical to the result obtained from synchronous experiments.

281 • **tabulate**: interpret the results of a benchmark and summarize them  
282 in the CSV file.

283 We present usage of each of the above commands later in section 4.

#### 284 3.2.4. *Asynchronous vs. synchronous execution*

285 PyQBench’s CLI can be used in synchronous and asynchronous modes.  
286 The mode of execution is defined in the YAML file describing the backend  
287 (see Section 4 for an example of this configuration). We decided to couple  
288 the mode of execution to the backend description because some backends  
289 cannot work in asynchronous mode.

290 When running `qbench disc-fourier benchmark` in asynchronous mode,  
291 the PyQBench submits all the circuits needed to perform a benchmark and  
292 then writes an intermediate YAML file containing metadata of submitted  
293 experiments. In particular, this metadata contains information on correlating  
294 submitted job identifiers with particular circuits. The intermediate file can  
295 be used to query the status of the submitted jobs or to resolve them, i.e. to  
296 wait for their completion and get the measurement outcomes.

297 In synchronous mode, PyQBench first submits all jobs required to run the  
298 benchmark and then immediately waits for their completion. The advantage  
299 of this approach is that no separate invocation of `qbench` command is needed  
300 to actually download the measurement outcomes. The downside, however,  
301 is that if the script is interrupted while the command is running, the inter-  
302 mediate results will be lost. Therefore, we recommend using asynchronous  
303 mode whenever possible.

## 304 4. Illustrative examples

305 In this section, we demonstrate the usage of PyQBench. For brevity, we  
306 decided to present only the usage of the CLI tool, as it is likely to be the  
307 most popular use case. We refer readers interested in implementing their

discrimination schemes using custom measurements to PyQBench’s documentation [43], where we describe the whole process, and to the Section 3 in the supplemental material, where we discuss the relevant mathematical details.

#### 4.1. Using *qbench* CLI

~~PyQBench offers a simplified way of conducting benchmarks using a Command-Line Interface (CLI).~~ The workflow with PyQBench’s CLI can be summarized as the following list of steps:

1. Preparing configuration files describing the backend and the experiment scenario.
2. Submitting/running experiments. Depending on the experiment scenario, execution can be synchronous, or asynchronous.
3. (optional) Checking the status of the submitted jobs if the execution is asynchronous.
4. Resolving asynchronous jobs into the actual measurement outcomes.
5. Converting obtained measurement outcomes into tabulated form.

##### 4.1.1. Preparing configuration files

The configuration of PyQBench CLI is driven by YAML files. The first configuration file describes the experiment ~~scenario to be executed~~scenario. The second file describes the backend. Typically, this backend will correspond to the physical device to be benchmarked, but for testing purposes, one might as well use any other Qiskit-compatible backend including simulators. ~~Let us first describe the~~The experiment configuration file, which might look as ~~follow~~follows.

Listing 2: Defining the experiment

---

```

type: discrimination-fourier
qubits:
  - target: 0
    ancilla: 1
  - target: 1
    ancilla: 2
angles:
  start: 0
  stop: 2 * pi
  num_steps: 3
gateset: ibmq
method: direct_sum
num_shots: 100

```

---

347 The second configuration file describes the backend. We decided to de-  
348 couple the experiment and the backend files because it facilitates their reuse.  
349 ~~For instance, the~~ The same experiment file can be used to run benchmarks  
350 on multiple backends, and the same backend description file can be used with  
351 multiple experiments.

352 Different Qiskit backends typically require different data for their initial-  
353 ization. Hence, there are multiple possible formats of the backend config-  
354 uration files understood by PyQBench. We refer the interested reader to  
355 the PyQBench's documentation. Below we describe an example YAML file  
356 describing IBM Q backend named Quito.

Listing 3: IBMQ backend

```
357 

---


358 name: ibmq_quito
359 asynchronous: false
360 provider:
361   hub: ibm-q
362   group: open
363   project: main
364 

---


```

365 IBMQ backends typically require an access token to IBM Quantum Experi-  
366 ence. Since it would be unsafe to store it in plain text, the token has to be  
367 configured separately in `IBMQ_TOKEN` environmental variable.

#### 368 *4.1.2. Running the experiment and collecting measurements data*

369 After preparing YAML files defining an experiment and backend, run-  
370 ning the benchmark can be launched by using the following command line  
371 invocation:

```
372 

---


373 qbench disc-fourier benchmark experiment_file.yml backend_file.yml
374 

---


```

375 The output file will be printed to stdout. Optionally, the `- -output OUTPUT`  
376 parameter might be provided to write the output to the `OUTPUT` file instead.

```
377 

---


378 qbench disc-fourier benchmark experiment_file.yml backend_file.yml
379   --output async_results.yml
380 

---


```

381 The result of running the above command can be twofold:

- 382 • If the backend is asynchronous, the output will contain intermediate  
383 data containing, amongst others, `job_ids` correlated with the circuit  
384 they correspond to.
- 385 • If the backend is synchronous, the output will contain measurement  
386 outcomes (bitstrings) for each of the circuits run.

For the synchronous experiment, the part of the output looks similar to the one below. The whole YAML file can be seen in Section ?? in the supplemental material.

---

```

data:
- target: 0
  ancilla: 1
  phi: 0.0
  results_per_circuit:
  - name: id
    histogram: {'00': 28, '01': 26, '10': 21, '11': 25}
    mitigation_info:
      target: {prob_meas0_prep1: 0.052200000000000024,
        prob_meas1_prep0: 0.0172}
      ancilla: {prob_meas0_prep1: 0.059000000000000005,
        prob_meas1_prep0: 0.0202}
    mitigated_histogram: {'00': 0.2637212373658018, '01':
      0.25865061319892463, '10': 0.2067279352110304, '11':
      0.2709002142242433}

```

---

#### 4.1.3. (Optional) Getting status of asynchronous jobs

PyQBench provides also a helper command that will fetch the statuses of asynchronous jobs. The command is:

---

```

qbench disc-fourier status async_results.yml

```

---

and it will display dictionary with histogram of statuses.

#### 4.1.4. Resolving asynchronous jobs

For asynchronous experiments, the stored intermediate data has to be resolved in actual measurements' outcomes. The following command will wait until all jobs are completed and then write a result file.

---

```

qbench disc-fourier resolve async-results.yml resolved.yml

```

---

The resolved results, stored in `resolved.yml`, would look just like if the experiment was run synchronously. Therefore, the final results will look the same no matter in which mode the benchmark was run, and hence in both cases, the final output file is suitable for being an input for the command computing the discrimination probabilities.

#### 4.1.5. Computing probabilities

As a last step in the processing workflow, the results file has to be passed to `tabulate` command:

---

```
qbench disc-fourier tabulate results.yml results.csv
```

---

A sample CSV file is provided in Table 2.

target	ancilla	phi	ideal_prob	disc_prob	mit_disc_prob
0	1	0	0.5	0.46	0.45
0	1	3.14	1	0.95	0.98
0	1	6.28	0.5	0.57	0.58
1	2	0	0.5	0.57	0.57
1	2	3.14	1	0.88	0.94
1	2	6.28	0.5	0.55	0.56

Table 2: The resulting CSV file contains table with columns `target`, `ancilla`, `phi`, `ideal_prob`, `disc_prob` and, optionally, `mit_disc_prob`. Each row in the table describes results for a tuple of (`target`, `ancilla`, `phi`). The reference optimal value of discrimination probability is present in `ideal_prob` column, whereas the obtained, empirical discrimination probability can be found in the `disc_prob` column. The `mit_disc_prob` column contains empirical discrimination probability after applying the `Mthree` error mitigation [46, 47], if it was applied.

432

## 5. Impact

With the surge of availability of quantum computing architectures in recent years it becomes increasingly difficult to keep track of their relative performance. ~~To make this case even more difficult~~ Moreover, various providers give access to different figures of merit for their architectures. Our package allows the user to test various architectures, available through `qiskit` and Amazon BraKet using problems with simple operational interpretation. We provide one example built-in in the package. Furthermore, we provide a powerful tool for the users to extend the range of available problems in a way that suits their needs.

Due to this possibility of extension, the users ~~are able to~~ can test specific aspects of their architecture of interest. For example, if their problem is related to the amount of coherence (the sum of absolute value of off-diagonal elements) of the states present during computation, they are able to quickly prepare a custom experiment, launch it on desired architectures, gather the



448 result, based on which they can decide which specific architecture they should  
449 use.

450 Finally, we provide the source code of PyQBench on GitHub [42] under an  
451 open source license which will allow users to utilize and extend our package  
452 in their specific applications.

## 453 6. Conclusions

454 In this paper, we presented a Python library PyQBench, an innovative  
455 open-source framework for benchmarking gate-based quantum computers.  
456 PyQBench can benchmark NISQ devices by verifying their capability of dis-  
457 criminating between two von Neumann measurements. PyQBench offers a  
458 simplified, ready-to-use, command line interface (CLI) for running bench-  
459 marks using a predefined parameterized Fourier family of measurements. For  
460 more advanced scenarios, PyQBench offers a way of employing user-defined  
461 measurements instead of predefined ones.

## 462 7. Conflict of Interest

463 We wish to confirm that there are no known conflicts of interest associated  
464 with this publication and there has been no significant financial support for  
465 this work that could have influenced its out- come.

## 466 Acknowledgements

467 This work is supported by the project “Near-term quantum computers  
468 Challenges, optimal implementations and applications” under Grant Num-  
469 ber POIR.04.04.00-00-17C1/18-00, which is carried out within the Team-Net  
470 programme of the Foundation for Polish Science co-financed by the Euro-  
471 pean Union under the European Regional Development Fund. PL is also a  
472 holder of European Union scholarship through the European Social Fund,  
473 grant InterPOWER (POWR.03.05.00-00-Z305).

## 474 References

- 475 [1] John Preskill. Quantum computing in the nisc era and beyond. *Quan-*  
476 *tum*, 2:79, 2018.
- 477 [2] Rigetti computing. <https://www.rigetti.com/>. Accessed on 2023-02-  
478 18.
- 479 [3] IBM Quantum. <https://www.ibm.com/quantum>. Accessed on 2023-02-  
480 18.

- 481 [4] Oxford quantum. <http://oxfordquantum.org/>. Accessed on 2023-02-  
482 18.
- 483 [5] IonQ. <https://ionq.com/>. Accessed on 2023-02-18.
- 484 [6] Xanadu. <https://www.xanadu.ai/>. Accessed on 2023-02-18.
- 485 [7] D-Wave Systems. <https://www.dwavesys.com/>. Accessed on 2023-02-  
486 18.
- 487 [8] QuEra. <https://www.quera.com/>. Accessed on 2023-02-18.
- 488 [9] QCS Documentation. <https://docs.rigetti.com/qcs/>. Accessed on  
489 2023-02-18.
- 490 [10] PyQuil Documentation. [https://pyquil-docs.rigetti.com/en/  
491 stable/](https://pyquil-docs.rigetti.com/en/stable/). Accessed on 2023-02-18.
- 492 [11] Qiskit. <https://qiskit.org/>. Accessed on 2023-02-18.
- 493 [12] IBM Quantum Experience. <https://quantum-computing.ibm.com/>.  
494 Accessed on 2023-02-18.
- 495 [13] Amazon Braket. <https://aws.amazon.com/braket/>. Accessed on  
496 2023-02-18.
- 497 [14] Zapata Orquestra Platform. [https://www.zapatacomputing.com/  
498 orquestra-platform/](https://www.zapatacomputing.com/orquestra-platform/). Accessed on 2023-02-18.
- 499 [15] Alexander J McCaskey, Dmitry I Lyakh, Eugene F Dumitrescu, Sarah S  
500 Powers, and Travis S Humble. XACC: a system-level software infrastruc-  
501 ture for heterogeneous quantum–classical computing. *Quantum Science  
502 and Technology*, 5(2):024002.
- 503 [16] The CUDA Quantum development team. CUDA Quantum. [https:  
504 //github.com/NVIDIA/cuda-quantum](https://github.com/NVIDIA/cuda-quantum). Accessed on 2023-08-27.
- 505 [17] J. Preskill. Quantum computing 40 years later. *arXiv preprint  
506 arXiv:2106.10522*, 2021.
- 507 [18] Y. Liu, M. Otten, R. Bassirianjahromi, L. Jiang, and B. Fefferman.  
508 Benchmarking near-term quantum computers via random circuit sam-  
509 pling. *arXiv preprint arXiv:2105.05232*, 2021.

- [19] E. Knill, D. Leibfried, R. Reichle, J. Britton, R. B. Blakestad, J. D. Jost, C. Langer, R. Ozeri, S. Seidelin, and D. J. Wineland. Randomized benchmarking of quantum gates. *Physical Review A*, 77(1):012307, 2008.
- [20] J. J. Wallman and S. T. Flammia. Randomized benchmarking with confidence. *New Journal of Physics*, 16(10):103032, 2014.
- [21] J. Helsen, I. Roth, E. Onorati, A. H. Werner, and J. Eisert. General framework for randomized benchmarking. *PRX Quantum*, 3(2):020357, 2022.
- [22] A. Cornelissen, J. Bausch, and A. Gilyén. Scalable benchmarks for gate-based quantum computers. *arXiv preprint arXiv:2104.10698*, 2021.
- [23] A. W. Cross, L. S. Bishop, S. Sheldon, P. D. Nation, and J. M. Gambetta. Validating quantum computers using randomized model circuits. *Physical Review A*, 100(3):032328, 2019.
- [24] N. Moll, P. Barkoutsos, L. S. Bishop, J. M. Chow, A. Cross, D. J. Egger, S. Filipp, A. Fuhrer, J. M. Gambetta, M. Ganzhorn, A. Kandala, A. Mezzacapo, P. Müller, W. Riess, G. Salis, J. Smolin, I. Tavernelli, and K. Temme. Quantum optimization using variational algorithms on near-term quantum devices. *Quantum Science and Technology*, 3(3):030503, 2018.
- [25] E. Pelofske, A. Bäertschi, and S. Eidenbenz. Quantum volume in practice: What users can expect from nist devices. *IEEE Transactions on Quantum Engineering*, 3:1–19, 2022.
- [26] A Yu Chernyavskiy and Yu I Bogdanov. Quantum tomography benchmarking. *Quantum Information Processing*, 20:1–20, 2021.
- [27] N. Quetschlich, L. Burgholzer, and R. Wille. Mqt bench: Benchmarking software and design automation tools for quantum computing. *arXiv preprint arXiv:2204.13719*, 2022.
- [28] MQTBench. <https://github.com/cda-tum/MQTBench>. Accessed on 2023-02-18.
- [29] T. Tomesh, P. Gokhale, V. Omole, G. S. Ravi, K. N. Smith, J. Visslai, X. Wu, N. Hardavellas, M. R. Martonosi, and F. T. Chong. SupermarQ: A scalable quantum benchmark suite. In *2022 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, pages 587–603. IEEE, 2022.

- [30] SupermarQ. <https://github.com/SupertechLabs/SupermarQ>. Accessed on 2023-02-18.
- [31] Qiskit benchmarks. <https://github.com/qiskit-community/qiskit-benchmarks>. Accessed on 2023-02-18.
- [32] Forest Benchmarking: QCVV using PyQuil. <https://github.com/rigetti/forest-benchmarking>. Accessed on 2023-02-18.
- [33] Quantum tomography benchmarking. <https://github.com/PQCLab/pyQTB>. Accessed on 2023-08-29.
- [34] Tirthak Patel, Abhay Potharaju, Baolin Li, Rohan Basu Roy, and Devesh Tiwari. Experimental evaluation of nisq quantum computers: Error measurement, characterization, and implications. In *SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–15. IEEE, 2020.
- [35] Ang Li, Samuel Stein, Sriram Krishnamoorthy, and James Ang. Qasm-bench: A low-level quantum benchmark suite for nisq evaluation and simulation. *ACM Transactions on Quantum Computing*, 4(2):1–26, 2023.
- [36] Quantum volume in practice. <https://github.com/lanl/Quantum-Volume-in-Practice>. Accessed on 2023-02-18.
- [37] Sergio Boixo, Sergei V Isakov, Vadim N Smelyanskiy, Ryan Babbush, Nan Ding, Zhang Jiang, Michael J Bremner, John M Martinis, and Hartmut Neven. Characterizing quantum supremacy in near-term devices. *Nature Physics*, 14(6):595–600, 2018.
- [38] Frank Arute, Kunal Arya, Ryan Babbush, Dave Bacon, Joseph C Bardin, Rami Barends, Rupak Biswas, Sergio Boixo, Fernando GSL Brandao, David A Buell, et al. Quantum supremacy using a programmable superconducting processor. *Nature*, 574(7779):505–510, 2019.
- [39] Colin Kai-Uwe Becker, Nikolay Tcholtchev, Ilie-Daniel Gheorghe-Pop, Sebastian Bock, Raphael Seidel, and Manfred Hauswirth. Towards a quantum benchmark suite with standardized kpis. In *2022 IEEE 19th International Conference on Software Architecture Companion (ICSAC)*, pages 160–163, 2022.

- 577 [40] Z. Puchała, Ł. Paweł, A. Krawiec, and R. Kukulski. Strategies for  
578 optimal single-shot discrimination of quantum measurements. *Physical*  
579 *Review A*, 98(4):042103, 2018.
- 580 [41] YAML Ain't Markup Language (YAML) Version 1.2. [https://yaml.](https://yaml.org/spec/1.2.2/)  
581 [org/spec/1.2.2/](https://yaml.org/spec/1.2.2/). Accessed on 2023-02-18.
- 582 [42] PyQBench GitHub Repository. [https://github.com/iitis/](https://github.com/iitis/PyQBench)  
583 [PyQBench](https://github.com/iitis/PyQBench). Accessed on 2023-02-18.
- 584 [43] PyQBench Documentation. [https://pyqbench.readthedocs.io/en/](https://pyqbench.readthedocs.io/en/latest/)  
585 [latest/](https://pyqbench.readthedocs.io/en/latest/). Accessed on 2023-02-18.
- 586 [44] Introducing the Qiskit Provider for Amazon Braket.  
587 [https://aws.amazon.com/blogs/quantum-computing/](https://aws.amazon.com/blogs/quantum-computing/introducing-the-qiskit-provider-for-amazon-braket/)  
588 [introducing-the-qiskit-provider-for-amazon-braket/](https://aws.amazon.com/blogs/quantum-computing/introducing-the-qiskit-provider-for-amazon-braket/). Ac-  
589 cessed on 2023-02-18.
- 590 [45] Qiskit Braket Provider GitHub Repository. [https://github.com/](https://github.com/qiskit-community/qiskit-braket-provider)  
591 [qiskit-community/qiskit-braket-provider](https://github.com/qiskit-community/qiskit-braket-provider). Accessed on 2023-02-  
592 18.
- 593 [46] mthree Documentation. [https://qiskit.org/documentation/](https://qiskit.org/documentation/partners/mthree/stubs/mthree.M3Mitigation.html)  
594 [partners/mthree/stubs/mthree.M3Mitigation.html](https://qiskit.org/documentation/partners/mthree/stubs/mthree.M3Mitigation.html). Accessed on  
595 2023-02-10.
- 596 [47] P. D. Nation, H. Kang, N. Sundaresan, and J. M. Gambetta. Scalable  
597 mitigation of measurement errors on quantum computers. *PRX Quan-*  
598 *tum*, 2(4):040326, 2021.