



Solver Computation Time

USER MANUAL

2020-09-25

Overview

This document provides timing information for D-Wave solvers. For QPU solvers it describes the timing of computation steps, including the overall service time allocated to problems, and how the QPU is accessed within that period. It lists timing-related parameters of QPU and Leap's hybrid solvers, including those that allow for user control.

CONTACT

Corporate Headquarters
3033 Beta Ave
Burnaby, BC V5G 4M9
Canada
Tel. 604-630-1428

US Office
2650 E Bayshore Rd
Palo Alto, CA 94303

Email: info@dwavesys.com

www.dwavesys.com

Notice and Disclaimer

D-Wave Systems Inc. (D-Wave), its subsidiaries and affiliates, makes commercially reasonable efforts to ensure that the information in this document is accurate and up to date, but errors may occur. NONE OF D-WAVE SYSTEMS INC., its subsidiaries and affiliates, OR ANY OF ITS RESPECTIVE DIRECTORS, EMPLOYEES, AGENTS, OR OTHER REPRESENTATIVES WILL BE LIABLE FOR DAMAGES, CLAIMS, EXPENSES OR OTHER COSTS (INCLUDING WITHOUT LIMITATION LEGAL FEES) ARISING OUT OF OR IN CONNECTION WITH THE USE OF THIS DOCUMENT OR ANY INFORMATION CONTAINED OR REFERRED TO IN IT. THIS IS A COMPREHENSIVE LIMITATION OF LIABILITY THAT APPLIES TO ALL DAMAGES OF ANY KIND, INCLUDING (WITHOUT LIMITATION) COMPENSATORY, DIRECT, INDIRECT, EXEMPLARY, PUNITIVE AND CONSEQUENTIAL DAMAGES, LOSS OF PROGRAMS OR DATA, INCOME OR PROFIT, LOSS OR DAMAGE TO PROPERTY, AND CLAIMS OF THIRD PARTIES.

D-Wave reserves the right to alter this document and other referenced documents without notice from time to time and at its sole discretion. D-Wave reserves its intellectual property rights in and to this document and its proprietary technology, including copyright, trademark rights, industrial design rights, and patent rights. D-Wave trademarks used herein include D-Wave®, Leap™ quantum cloud service, Ocean™, Advantage™ quantum system, D-Wave 2000Q™, D-Wave 2X™, and the D-Wave logos (the D-Wave Marks). Other marks used in this document are the property of their respective owners. D-Wave does not grant any license, assignment, or other grant of interest in or to the copyright of this document, the D-Wave Marks, any other marks used in this document, or any other intellectual property rights used or referred to herein, except as D-Wave may expressly provide in a written agreement. This document may refer to other documents, including documents subject to the rights of third parties. Nothing in this document constitutes a grant by D-Wave of any license, assignment, or any other interest in the copyright or other intellectual property rights of such other documents. Any use of such other documents is subject to the rights of D-Wave and/or any applicable third parties in those documents.

Contents

1	About this Document	1
1.1	Intended Audience	1
1.2	Scope	1
1.3	Glossary	1
2	QMI Timing Overview	3
3	Breakdown of QPU Access Time	5
4	Sources of Timing Variation and Error	7
4.1	Nondedicated QPU Use	7
4.2	Nondeterminacy of Classical System Timings	7
4.3	Internet Latency	9
4.4	Settings of User-Specified Parameters	9
5	Breakdown of Service Time	11
5.1	Postprocessing Time	11
5.2	“Total Time” Reported in Statistics (for Administrators)	12
6	QPU Timing Information from Ocean Software	15
6.1	Timing-Related Fields	15
6.2	Timing Data Returned by <code>dwave-cloud-client</code>	17
7	Leap’s Hybrid Solvers’ Timing Information	19
7.1	Timing-Related Fields	19

ABOUT THIS DOCUMENT

1.1 Intended Audience

This document is for users of D-Wave™ solvers who want to understand how their use of the system is timed, which bears on accounting.

1.2 Scope

This document describes the computation process of D-Wave quantum computers, focusing on system timing. It explains the overall time that is allocated to a quantum machine instruction (QMI), describes how use of the quantum processing unit (QPU) is timed within that period, gives context for how timing can vary.

For QPU and Leap's hybrid solvers it describes the timing-related fields in the Solver API (SAPI).

1.3 Glossary

This document uses the following terms:

- **Hybrid solvers:** Quantum-classical hybrid solvers implement state-of-the-art classical algorithms together with intelligent allocation of the QPU to parts of the problem where it benefits most.
- **Quantum processing unit (QPU):** Quantum computational element within a D-Wave system.
- **Quantum machine instruction (QMI):** Set of information that is sent to the QPU, including the Ising model or QUBO parameters and annealing-cycle parameters.
- **Annealing cycle:** Physical process of starting with the QMI (prepared for the D-Wave system) and yielding a single sample to the QMI. Typically executed multiple times in a single QMI.
- **Sample, result, read, or solution:** Terms for the result yielded by an annealing cycle; 'sample' is preferred to connote the nondeterministic behavior of the QPU.

QMI TIMING OVERVIEW

Fig. 2.1 shows a simplified diagram of the sequence of steps, the dark red set of arrows, to execute a quantum machine instruction (QMI) on a D-Wave system, starting and ending on a user's client system. Each QMI consists of a single input together with parameters. A QMI is sent across a network to the SAPI server and joins a queue. Each queued QMI is assigned to one of possibly multiple workers, which may run in parallel. A worker prepares the QMI for the quantum processing unit (QPU) and optionally for postprocessing¹, sends the QMI to the QPU queue, receives samples (results) and optionally post-processes them (overlapping in time with QPU execution), and bundles the samples with additional QMI-execution information for return to the client system.

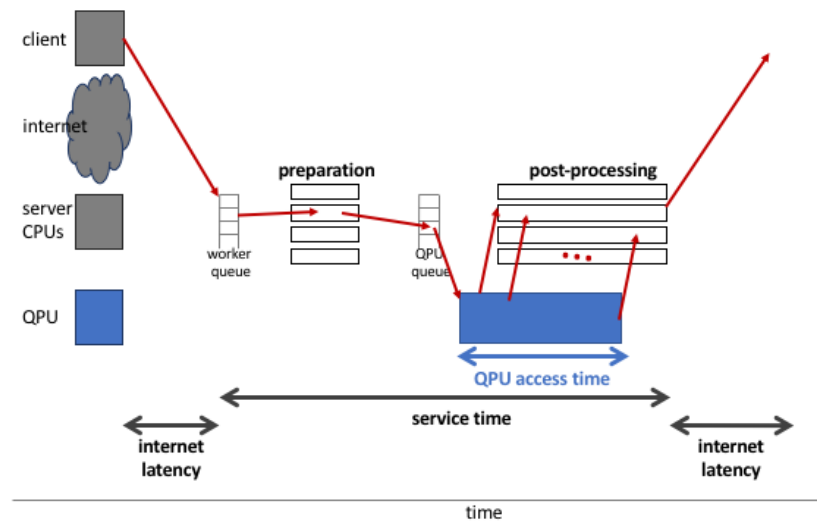


Figure 2.1: Overview of execution of a single QMI, starting from a client system, and distinguishing classical (client, CPU) and quantum (QPU) execution.

The total time for a QMI to pass through the D-Wave system is the *service time*. The execution time for a QMI as observed by a client includes service time and internet latency. The QPU executes one QMI at a time, during which the QPU is unavailable to any other QMI. This execution time is known as the QMI's *QPU access time*.

¹ Postprocessing for D-Wave 2000Q and earlier systems includes optimization and sampling algorithms; on Advantage systems, postprocessing is limited to computing the energies of returned samples. Ocean software provides postprocessing tools, you can use these for Advantage systems.

Important: Users executing on a D-Wave system are charged for QPU access time.

For customers purchasing time from D-Wave, this will be a financial charge. For users on a customer-owned system, this will typically be an administrative charge.

BREAKDOWN OF QPU ACCESS TIME

As illustrated in Figure 3.1, the time to execute a single QMI on a QPU, *QPU access time*, is broken into two parts: a one-time initialization step to program the QPU (blue), and typically multiple (R) per-sample times for the actual execution on the QPU (repeated multi-color).

Note: A small amount of initialization time spent in low-level operations, denoted as Δ in the equations below, is roughly 1 ms and is reported as *qpu_access_overhead_time* by SAPI.

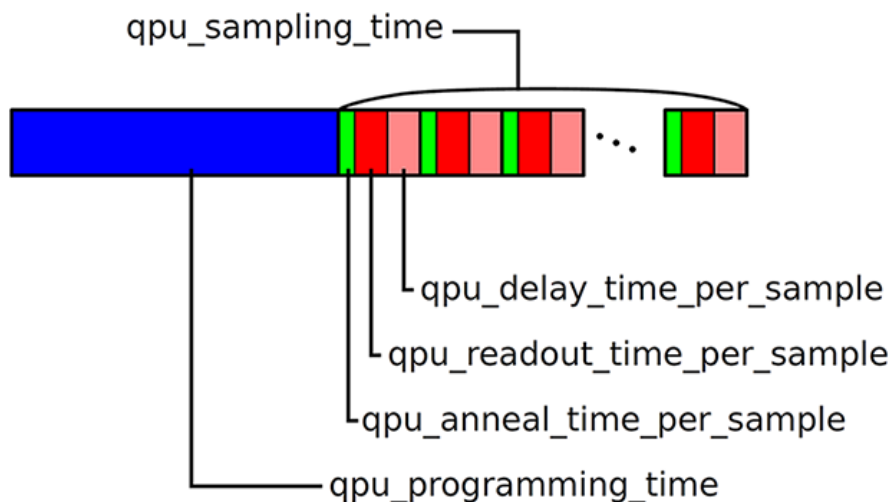


Figure 3.1: Detail of QPU access time.

The end-to-end time for one sample is further broken into anneal (the anneal proper; green), readout (read the sample from the QPU; red), and thermalization (wait for the QPU to regain its initial temperature; pink). Note that these component times are calculated as average times per sample. Possible rounding errors mean that the sum of these times may not match the total sampling time reported. Thus, for total time T , these times relate as follows:

$$T = T_p + \Delta + T_s$$

$$T_s/R \approx T_a + T_r + T_d,$$

where

$$\begin{aligned} T_p &= \text{programming_time} \\ T_s &= \text{sampling_time} \end{aligned} \tag{3.1}$$

and

$$\begin{aligned} T_a &= \text{(total) annealing_time} \\ T_r &= \text{(total) readout_time} \\ T_d &= \begin{aligned} &\text{(total) anneal_schedule} \\ &+ \text{(total) readout_thermalization} \\ &+ \text{(total) reduce_intersample_correlation} \\ &+ \text{(total) reinitialize_state.} \end{aligned} \end{aligned} \tag{3.2}$$

Note that *reinitialize_state* is used only for reverse annealing.

SOURCES OF TIMING VARIATION AND ERROR

Running a D-Wave-using program across the internet or even examining QPU timing information may show variation from run to run from the end-user's point of view. This section describes some of the possible sources of such variation.

4.1 Nondedicated QPU Use

D-Wave systems are typically shared among multiple users, each of whom submits QMIs to solve a problem, with little to no synchronization among users. (A single user may also have multiple client programs submitting unsynchronized QMIs to a D-Wave system.) The QPU must be used by a single QMI at a time, so the D-Wave system software ensures that multiple QMIs flow through the system and use the QPU sequentially. In general, this means that a QMI may get queued for the QPU or some other resource, injecting indeterminacy into the timing of execution.

Note: Contact your D-Wave system administrator or D-Wave Support if you need to ensure a quiet system.

4.2 Nondeterminacy of Classical System Timings

Even when a system is quiet except for the program to be measured, timings often vary. As illustrated in Fig. 4.1, running a given code block repeatedly can yield different runtimes on a classical system, even though the instruction execution sequence does not change. Runtime distributions with occasional large outliers, as seen here, are not unusual.

Timing variations are routine, caused by noise from the operating system (e.g., scheduling, memory management, and power management) and the runtime environment (e.g., garbage collection, just-in-time compilation, and thread migration).¹ In addition, the internal architecture of the classical portion of the D-Wave system includes multiple hardware nodes and software servers, introducing communication among these servers as another source of variation.

¹ A more common practice in computational research is to report an alternative measurement called CPU time, which is intended to filter out operating system noise. However, CPU timers are only accurate to tens of milliseconds, and CPU times are not available for QPU time measurements. For consistency, we use wall clock times throughout.

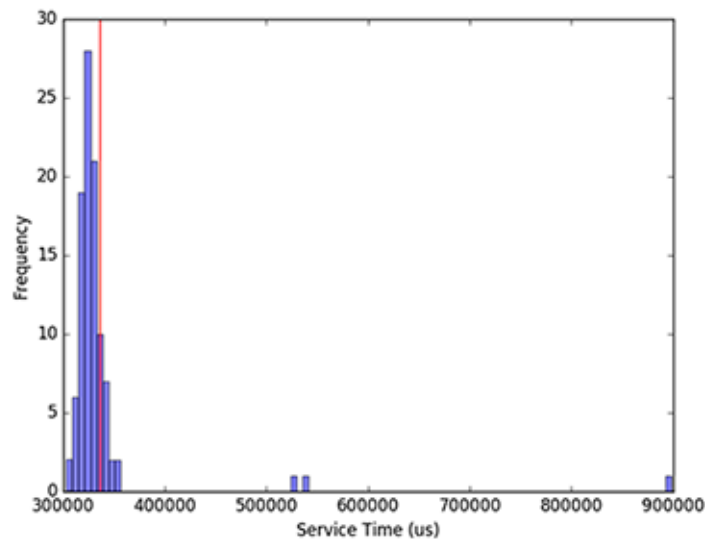


Figure 4.1: Histogram of 100 measurements of classical execution time using a wall clock timer, showing that the mean time of 336.5 ms (red line) is higher than 75 percent of the measurements.

For these reasons, mean reported runtimes can often be higher than median runtimes: for example, in Fig. 4.1, the mean time of 336.5 ms (vertical red line) is higher than 75 percent of the measured runtimes due to a few extreme outliers (one about 3 times higher and two almost 2 times higher than median). As a result, mean runtimes tend to exceed median runtimes. In this context, the smallest time recorded for a single process is considered the most accurate, because noise from outside sources can only increase elapsed time.² Because system activity increases with the number of active QMIs, the most accurate times for a single process are obtained by measuring on an otherwise quiet system.

Note: The 336 ms mean time shown for this particular QMI is not intended to be representative of QMI execution times.

The cost of reading a system timer may impose additional measurement errors, since querying the system clock can take microseconds. To reduce the impact of timing code itself, a given code block may be measured outside a loop that executes it many times, with running time calculated as the average time per iteration. Because of system and runtime noise and timer latency, component times measured one way may not add up to total times measured another way.³ These sources of timer variation or error are present on all computer systems, including the classical portion of D-Wave platforms. Normal timer variation as described here may occasionally yield atypical and imprecise results; also, one expects wall clock times to vary with the particular system configuration and with system load.

² Randal E. Bryant and David R. O'Hallaron, *Computer Systems: A Programmer's Perspective (2nd Edition)*, Pearson, 2010.

³ Paulo Eduardo Nogueira, Rivalino Matias, Jr., and Elder Vicente, *An Experimental Study on Execution Time Variation in Computer Experiments*, ACM Symposium on Applied Computing, 2014.

4.3 Internet Latency

If you are running your program on a client system geographically remote from the D-Wave system on which you're executing, you will likely encounter latency and variability from the internet connection itself (see Fig. 2.1).

4.4 Settings of User-Specified Parameters

The following user-specified parameters can cause timing to change, but should not affect the variability of timing. For more information on these parameters, see *Solver Properties and Parameters Reference*.

- *anneal_schedule*—User-provided anneal schedule. Specifies the points at which to change the default schedule. Each point is a pair of values representing time t in microseconds and normalized anneal fraction s . The system connects the points with piecewise-linear ramps to construct the new schedule. If *anneal_schedule* is specified, T_a , *qpu_anneal_time_per_sample* is populated with the total time specified by the piecewise-linear schedule. Cannot be used with the *annealing_time* parameter.
- *annealing_time*—Duration, in microseconds, of quantum annealing time. This value populates T_a , *qpu_anneal_time_per_sample*. Cannot be used with the *anneal_schedule* parameter.
- *num_reads*—Number of samples to read from the solver per QMI.
- *num_spin_reversal_transforms*—For QMIs with more than one spin-reversal transform, SAPI handles the timing information for all the subQMIs that it sends to the solver as follows: (1) It sums each timing field that does not end with “per_sample.” (2) For others, it sends the value from the first subQMI. For example, the values for *qpu_access_time* are summed; those from *qpu_delay_time_per_sample* are not.
- *postprocess*—Specifies the type of (classical) postprocessing to be performed on the raw samples from the QPU. Requesting no postprocessing consumes the least time; either sampling or optimization postprocessing consumes more.
- *programming_thermalization*—Number of microseconds to wait after programming the QPU to allow it to cool; i.e., *post-programming thermalization time*. Values lower than the default accelerate solving at the expense of solution quality. This value contributes to the total T_p , *qpu_programming_time*.
- *readout_thermalization*—Number of microseconds to wait after each sample is read from the QPU to allow it to cool to base temperature; i.e., *post-readout thermalization time*. This optional value contributes to T_d , *qpu_delay_time_per_sample*.
- *reduce_intersample_correlation*—Used to reduce sample-to-sample correlations. When true, adds to T_d , *qpu_delay_time_per_sample*. Amount of time added increases linearly with increasing length of the anneal schedule.
- *reinitialize_state*—Used in reverse annealing. When true (the default setting), reinitializes the initial qubit state for every anneal-readout cycle, adding between 100 and 600 microseconds to T_d , *qpu_delay_time_per_sample*. When false, adds approximately

10 microseconds to T_d .⁴

Note: Depending on the parameters chosen for a QMI, QPU access time may be a large or small fraction of service time. E.g., a QMI requesting a single sample with short *annealing_time* would see programming time as a large fraction of service time and QPU access time as a small fraction.

4.4.1 Upper Limit on User-Specified Timing-Related Parameters

The D-Wave system limits your ability to submit a long-running QMI to prevent you from inadvertently monopolizing QPU time. This limit varies by system; check the *problem_run_duration_range* property for your solver.

The limit is calculated according to the following formula:

$$Duration = ((P_1 + P_2) * P_3) + P_4 \quad (4.1)$$

where P_1 , P_2 , P_3 , and P_4 are the values specified for the *annealing_time*, *read-out_thermalization*, *num_reads* (samples), and *programming_thermalization* parameters, respectively.

If you attempt to submit a QMI whose execution time would exceed the limit for your system, an error is returned showing the values in microseconds. For example:

```
ERROR: Upper limit on user-specified timing related parameters exceeded: 12010000 > 3000000
```

Note that it is possible to specify values that fall within the permitted ranges for each individual parameter, yet together cause the time to execute the QMI to surpass the limit.

⁴ Amount of time varies by system.

BREAKDOWN OF SERVICE TIME

The service time can be broken into:

- Any time required by the worker before and after QPU access
- Wait time in queues before and after QPU access
- QPU access time
- Postprocessing (PP) time

Service time is defined as the difference between time-in and time-out for each QMI, as shown in the table.

Keyword in SAPI	Meaning
<i>time_received</i>	When QMI arrives
<i>time_solved</i>	When bundled samples are available

As mentioned in the introduction, service time for a single QMI depends on the system load; that is, how many other QMIs are present at a given time. During periods of heavy load, wait time in the two queues may contribute to increased service times. D-Wave has no control over system load under normal operating conditions. Therefore, it is not possible to guarantee that service time targets can be met. Service time measurements described in other D-Wave documents are intended only to give a rough idea of the range of experience that might be found under varying conditions.

5.1 Postprocessing Time

Postprocessing optimization and sampling algorithms provide local improvements with minimal overhead to solutions obtained from the quantum processing unit (QPU).

[Ocean software](#) provides postprocessing tools, and you can optionally run postprocessing online on D-Wave 2000Q and earlier systems.

As shown in [Fig. 5.1](#), online postprocessing (red) works in parallel with sampling (blue), so that the computation times overlap except for postprocessing the last batch of samples. In this diagram, the time consumed by gathering small batches of samples are marked by vertical blue lines. Within execution of a QMI, gathering the current set of samples takes place concurrently with postprocessing for the previous set of samples (red boxes), which is applied to batches of samples as they are returned by the QPU. As illustrated by [Fig. 5.1](#), only the time for postprocessing the last set of samples (the rightmost red box) is not overlapped with sampling.

Postprocessing overhead is designed not to impose any delay to QPU access for the next

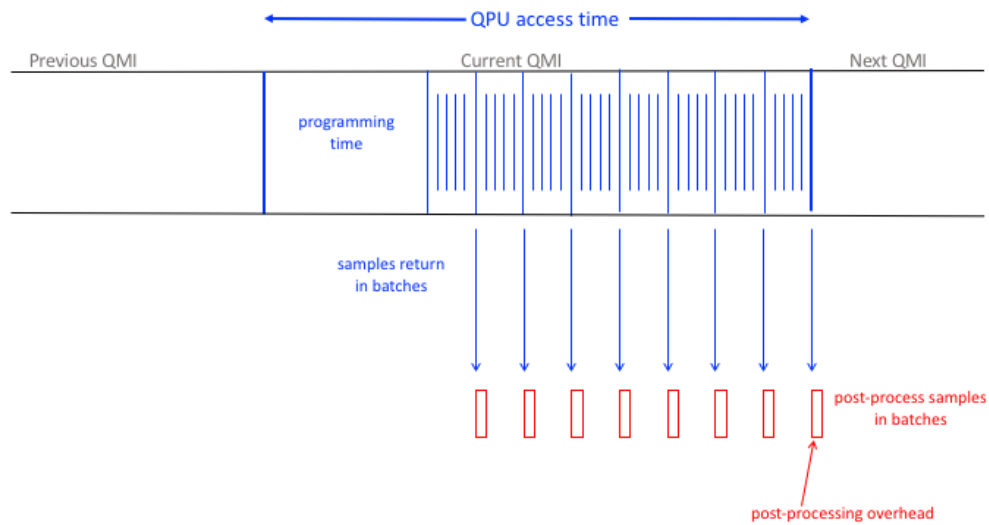


Figure 5.1: Relationship of QPU time to postprocessing time, illustrated by one QMI in a sequence (previous, current, next).

QMI, because postprocessing of the last batch takes place concurrently with the next QMI’s programming time.

The system returns two associated timing values, as shown in the table below. Referring to Fig. 5.1, *total_post_processing_time* is the sum of all times in the red boxes, while *post_processing_overhead* is the extra time needed (a single red box) to process the last batch. This latter time together with *qpu_access_time* contributes to overall service time.

Note: Even if no postprocessing is run on a QMI, the returned *post_processing_overhead* value is non-zero. This is because computing the final energies occurs after samples are returned and is accounted as postprocessing overhead.

Keyword in SAPI	Meaning
<i>total_post_processing_time</i>	Total time for postprocessing
<i>post_processing_overhead_time</i>	Added for the last batch

For more details about postprocessing and how it is handled in the timing structure, see *Postprocessing Methods on D-Wave Systems*.

5.2 “Total Time” Reported in Statistics (for Administrators)

One timing parameter, *qpu_access_time*, provides the raw data for the “Total Time” values reported as system statistics, available to administrators. Reported statistics are the sum of the *qpu_access_time* values for each QMI selected by the users, solvers, and time periods

selected in the filter.

Note: Reported statistics are in milliseconds, while SAPI inputs and outputs are in microseconds. One millisecond is 1000 microseconds.

QPU TIMING INFORMATION FROM OCEAN SOFTWARE

6.1 Timing-Related Fields

The table below lists the timing-related fields available in D-Wave's [Ocean SDK](#). Ocean users access to this information is from the `info` field in the `dimod` `sampleset` class, as in the example below.

```
>>> from dwave.system import DWaveSampler, EmbeddingComposite
>>> sampler = EmbeddingComposite(DWaveSampler(solver={'qpu': True}))
>>> sampleset = sampler.sample_ising({'a': 1}, {'a', 'b': 1})
>>> print(sampleset.info["timing"])
{'qpu_sampling_time': 315,
 'qpu_anneal_time_per_sample': 20,
 'qpu_readout_time_per_sample': 274,
 'qpu_access_time': 9687,
 'qpu_access_overhead_time': 1738,
 'qpu_programming_time': 9372,
 'qpu_delay_time_per_sample': 21,
 'total_post_processing_time': 474,
 'post_processing_overhead_time': 474,
 'total_real_time': 9687,
 'run_time_chip': 315,
 'anneal_time_per_run': 20,
 'readout_time_per_run': 274}
```

Table 6.1: Fields that affect *qpu_access_time*

Value	Field in dwave-cloud-client	Meaning	Affected by
T	<i>qpu_access_time</i>	Total time in QPU	All parameters listed below
T_p	<i>qpu_programming_time</i>	Total time to program the QPU ¹	<i>programming_thermalization</i>
Δ		Additional low-level operations	
R		Number of reads (samples)	
T_s	<i>qpu_sampling_time</i>	Total time for R samples	<i>num_reads</i>
T_a	<i>qpu_anneal_time_per_sample</i>	Time for one anneal	<i>anneal_schedule, annealing_time</i>
T_r	<i>qpu_readout_time_per_sample</i>	Time for one read	Number of qubits read ²
T_d	<i>qpu_delay_time_per_sample</i>	Delay between anneals ³	<i>anneal_schedule, readout_thermalization, reduce_intersample_correlation, reinitialize_state</i> (only in case of reverse annealing)

¹ Even if *programming_thermalization* is 0, T_p is several milliseconds for other operations.

² The time to read a sample set from a QPU depends on the number of qubits in the sample set: a problem represented by a dozen qubits has shorter read times (and so a shorter T_r , the *total_readout_time*) than a problem represented by several thousand qubits.

For the Advantage QPU, this can be significant. For example, some small problem might have reads of ~100 ms while for a large problem these might take ~300 ms.

³ The time returned in the *qpu_delay_time_per_sample* field is equal to a constant plus the user-specified value, *readout_thermalization*.

6.2 Timing Data Returned by `dwave-cloud-client`

Below is a sample skeleton of Python code for accessing timing data returned by `dwave-cloud-client`. Timing values are returned in the computation object and the timing object; further code could query those objects in more detail. The timing object referenced on line 16 is a Python dictionary containing (key, value) pairs. The keys match keywords discussed in this section.

```

01 import random
02 import datetime as dt
03 from dwave.cloud import Client

04 # Connect using the default or environment connection information
05 with Client.from_config() as client:

06     # Load the default solver
07     solver = client.get_solver()

08     # Build a random Ising model to exactly fit the graph the solver supports
09     linear = {index: random.choice([-1, 1]) for index in solver.nodes}
10     quad = {key: random.choice([-1, 1]) for key in solver.undirected_edges}

11     # Send the problem for sampling, include solver-specific parameter 'num_reads
    →,
12     computation = solver.sample_ising(linear, quad, num_reads=100)
13     computation.wait()

14     # Print the first sample out of a hundred
15     print(computation.samples[0])
16     timing = computation['timing']

17     # Service time
18     time_format = "%Y-%m-%d %H:%M:%S.%f"
19     start_time = dt.datetime.strptime(str(computation.time_received)[: -6], time_
    →format)
20     finish_time = dt.datetime.strptime(str(computation.time_solved)[: -6], time_
    →format)
21     service_time = finish_time - start_time
22     qpu_access_time = timing['qpu_access_time']
23     print("start_time="+str(start_time)+", finish_time="+str(finish_time)+ \
24           ", service_time="+str(service_time)+", qpu_access_time=" \
25           +str(float(qpu_access_time)/1000000))

```


LEAP'S HYBRID SOLVERS' TIMING INFORMATION

7.1 Timing-Related Fields

The table below lists the timing information returned from Leap's quantum-classical hybrid solvers. Ocean users access to this information is from the `info` field in the `dimod.sampleset` class, as in the example below.

```
>>> import dimod
>>> from dwave.system import LeapHybridSampler
...
>>> sampler = LeapHybridSampler(solver={'category': 'hybrid'})
>>> bqm = dimod.generators.ran_r(1, 300)
>>> sampleset = sampler.sample(bqm)
>>> sampleset.info
{'qpu_access_time': 41990, 'charge_time': 2991424, 'run_time': 2991424}
```

Table 7.1: Timing Information from Leap's Hybrid Solvers

Field	Meaning
<code>run_time</code>	Time, in microseconds, the hybrid solver spent working on the problem.
<code>charge_time</code>	Time, in microseconds, charged to the account. ¹
<code>qpu_access_time</code>	QPU time, in microseconds, used by the hybrid solver. ²

¹ `charge_time` and `run_time` may differ due to the time granularity of the solver. For example, if a hybrid solver has a time granularity of 0.5 sec and your submission specifies 4.7 sec, `run_time` might be 5 sec and `charge_time` 4.5 sec.

² `qpu_access_time` might be zero for submissions with short `run_time`. Because the QPU is a shared, remote resource, the first set of samples from the QPU might not be received before a short explicitly-specified or default `time_limit` is reached. In such cases, the hybrid solver respects the time limitation and returns without the QPU having a chance to contribute to the solution. On the large, complex problems for which hybrid solvers are intended, this is unlikely to occur.