# A Hybrid Quantum-Classical Paradigm to Mitigate Embedding Costs in Quantum Annealing—Abridged Version*

Alastair A. Abbott

Univ. Grenoble Alpes, CNRS, Grenoble INP, Institut Néel
38000 Grenoble, France

Cristian S. Calude     Michael J. Dinneen     Richard Hua

Department of Computer Science, University of Auckland
Private Bag 92019, Auckland, New Zealand

Quantum annealing has shown significant potential as an approach to near-term quantum computing. Despite promising progress towards obtaining a quantum speedup, quantum annealers are limited by the need to embed problem instances within the (often highly restricted) connectivity graph of the annealer. This embedding can be costly to perform and may destroy any computational speedup. Here we present a hybrid quantum-classical paradigm to help mitigate this limitation, and show how a raw speedup that is negated by the embedding time can nonetheless be exploited in certain circumstances. We illustrate this approach with initial results on a proof-of-concept implementation of an algorithm for the dynamically weighted maximum independent set problem.

## 1 Introduction

Quantum computation has the potential to revolutionise computer science, and as a consequence has received a great deal of attention from theorists and experimentalists alike. Although much progress has been made through the concerted efforts of the community, we are still some distance from being able to build sufficiently large-scale universal quantum computers to realise this potential [21].

More recently, however, significant progress has been made in the development of special-purpose quantum computers. This has been driven by the realisation that, by dropping the requirement of being able to efficiently simulate arbitrary computations and relaxing some of the constraints that make large-scale universal quantum computing difficult, such devices can be more easily engineered and scaled. With this approach it may be possible to exploit some of the capabilities of quantum computation to obtain lesser, but nevertheless practical, advantages in real-world applications. Quantum annealers, which solve particular optimisation problems, exemplify this approach, and significant progress has been made in recent years towards engineering moderately large-scale such devices [17]. This approach has been pursued particularly zealously by D-Wave, who have developed quantum annealers with upwards of 2000 qubits [11] and are thus of sufficient size to tackle problems for which their performance can meaningfully be compared to classical computational approaches.

In this paradigm, however, it is a subtle problem to compare the performance of quantum solutions to classical ones, since the focus is on obtaining real-world gains in domains where heuristics tend to be at the core of the best classical approaches. Indeed, this issue is at the heart of recent debate surrounding

---

*A significantly expanded version is available at arXiv:1803.04340 [cs.DS].

the performance of D-Wave machines [6, 27]. In particular, instead of focusing on asymptotic analyses, one must compare the performance of classical and quantum devices empirically. But performing benchmarking fairly is difficult, especially when there is often debate as to which classical algorithm should be taken for comparison [18, 20, 26].

In this paper, motivated by the need to take into account the cost of classical processing in benchmarking quantum annealers, we propose a hybrid quantum-classical approach for developing algorithms that mitigates the cost of this processing. We focus on D-Wave's quantum annealers where the process involves a costly classical "embedding" stage which maps an arbitrary problem instance into one compatible with D-Wave's limited connectivity constraints. We then formulate a generic hybrid approach that mitigates this cost allowing any advantage present to be accessed more directly [4]. The embedding problem is time-consuming, and experimental studies indicate that its quality can have strong effects on performance [30, 32].

To illustrate this generic framework for hybrid computing we propose, we present a hybrid algorithm based around a D-Wave solution to the maximum-weight independent set (MWIS) problem. We present an overview of the results of an initial proof-of-principle implementation of this algorithm, showing a large improvement of the hybrid algorithm over a more standard quantum annealing approach, as well as comparing it to a classical algorithm.

## 2   D-Wave's quantum annealing framework

### 2.1   Quantum annealing and quadratic unconstrained Boolean optimisation

Quantum annealing is a finite temperature implementation of adiabatic quantum computing [13], in which the optimisation problem to be solved is encoded into a Hamiltonian $H_p$ (the quantum operator corresponding to the system's energy) such that the ground state of $H_p$ corresponds precisely to the solution to the problem (or one of them, if many exist). The computer is initially prepared in the ground state of a Hamiltonian $H_i$, which is then slowly evolved into the target Hamiltonian $H_p$. This computation can be described by the time-dependent Hamiltonian $H(t) = A(t)H_i + B(t)H_p$ for $0 \leq t \leq T$, where $A(0) = B(T) = 1$ and $A(T) = B(0) = 0$. $T$ is called the annealing time and, for D-Wave machines, the functions $A$ and $B$ give a close to a linear transition from $H_i$ to $H_p$ [17]. If the computation is performed sufficiently slowly, the Adiabatic Theorem guarantees that the system will remain in a ground state of $H_p$ throughout the computation and the final state will thus correspond to an optimal solution to the problem at hand [13].

Quantum annealers implement specific, simple classes of Hamiltonians, such as the two-dimensional Ising spin Hamiltonians realised by D-Wave devices. This restriction means that D-Wave annealers are capable of solving natively the *Quadratic Unconstrained Boolean Optimisation (QUBO) problem* [7]. The QUBO problem is the task of finding the input $\mathbf{x}^*$ that minimises a quadratic objective function of the form $f(\mathbf{x}) = \mathbf{x}^T Q \mathbf{x}$, where $\mathbf{x} = (x_1, \ldots, x_n)$ is a vector of $n$ binary variables and $Q$ is an upper-triangular $n \times n$ matrix of real numbers:

$$\mathbf{x}^* = \arg\min_{\mathbf{x}} \mathbf{x}^T Q \mathbf{x} = \arg\min_{\mathbf{x}} \sum_{i \leq j} x_i Q_{(i,j)} x_j, \text{ where } x_i \in \{0, 1\}.$$

In the quantum annealing model of the QUBO problem, each $x_i$ corresponds to a qubit while $Q$ defines the problem Hamiltonian $H_p$. Crucially, the nonzero terms $Q_{(i,j)}$ (for $i \neq j$) correspond to couplings between qubits and induce a graph $G_L = (V_L, E_L)$ called the *logical graph* representing interactions between qubits; the qubits $V_L$ in which the QUBO problem is represented over are called the *logical qubits*.

## 2.2 Hardware constraints and embeddings

The comparative ease in engineering devices which naturally solve the QUBO problem has been crucial for the recent experimental success of quantum annealing. Still, it remains exceedingly difficult to control interactions between qubits that are not physically near to one another, and as a result it is not possible to implement directly any instance of the QUBO problem. Instead, the couplings possible on a quantum annealer are specified by a *physical graph* $G_P = (V_P, E_P)$, where $V_P$ is the set of *physical qubits* on the device, and an edge $\{i, j\} \in E_P$ signifies that qubits $i$ and $j$ can be physically coupled [7].

The physical graphs implemented on D-Wave's devices are *Chimera graphs* $\chi_k$, which are $k \times k$ grids of $K_{4,4}$ graphs [5]. Specifically, each qubit is coupled with four other qubits in the same $K_{4,4}$ block and two qubits in adjacent blocks (except for qubits in blocks on the edge of the grid, which are coupled to a single other block). The Chimera graph is crucially relatively sparse and near-to-planar, with qubits separated by paths of length up to $2k$. Since the logical graph $G_L$ for a QUBO problem instance $Q$ will not, in general, be a subgraph of the physical graph $G_P = \chi_k$, the problem instance on $G_L$ must be mapped to an equivalent one on $G_P$. This process involves two steps: first, $G_L$ must be *embedded* in $G_P$, and secondly the weights of the QUBO problem (i.e., the non-zero entries in $Q$) must be adjusted so that valid solutions on $G_P$ are mapped to valid solutions on $G_L$.

The embedding stage amounts to finding a *(minor) embedding* of $G_L$ into $G_P$ [7], i.e., an embedding function $f : V_L \to 2^{V_P}$ such that i) the sets of vertices $\{f(v) \mid v \in V_L\}$ are disjoint, ii) for all $v \in V_L$, there is a subset of edges $E' \subseteq E_P$ such that $G' = (f(v), E')$ is connected and iii) if $\{u, v\} \in E_L$, then there exist $u', v' \in V_P$ such that $u' \in f(u)$, $v' \in f(v)$ and $\{u', v'\}$ is an edge in $E_P$. The problem of finding a minor embedding is itself computationally difficult [7]. The embedding process may thus, in light of its computational difficulty, contribute significantly to the time required to solve a problem in practice. Currently, the standard approach to finding such an embedding is to use heuristic algorithms.

## 2.3 Benchmarking quantum annealers

It is not generally believed that an exponential speedup is possible for NP-hard problems such as the QUBO problem [1], and there has been much debate over whether or not quantum annealing provides any such speedup in practice [4, 26]. Indeed, there is disagreement over what exactly constitutes a quantum "speedup" and how to determine if there is one. In this paper we will focus primarily on the empirical run-time performance in investigating whether a quantum speedup is present, rather than the (empirically estimated) scaling performance of quantum algorithms.

Good benchmarking needs to make use of fair and comprehensive metrics to determine the running time of both classical and quantum algorithms for a problem. In particular these need to properly take into account not only the "wall-time" of different stages of the quantum algorithm, but also its probabilistic nature. To understand how this can be done, we first need to outline the different stages of the quantum annealing process [19]: 1) *Programming:* The problem instance is loaded onto the annealing chip (QPU),

which takes time $t_{\mathrm{prog}}$; 2) *Annealing:* The quantum annealing process is performed and then the physical qubits are measured to obtain a solution; this takes time $t_a$; 3) *Repetition:* Step 2 is performed $k$ times to obtain $k$ potential solutions. The *quantum processing time* (QPT) is thus $t_{\mathrm{proc}} = t_{\mathrm{prog}} + k t_a$. With these considerations on hand, a relatively fair and robust way to measure the quantum processing time is the "time to solution" (TTS) metric [3, 26], which is based on the expected number of repetitions needed to obtain a valid solution with probability $p$ (one typically takes $p = 0.99$). If the probability per annealing sample of obtaining a solution is $s$ (which can be estimated empirically), then this is calculated as $k_{99} = \frac{\log(1-p)}{\log(1-s)}$, and the quantum processing time is thus calculated with this $k$ as $t_{\mathrm{proc}} = t_{\mathrm{prog}} + k_{99} t_a$. Existing investigations have primarily focused on comparing directly the QPT with the processing time of a classical algorithm in order to look for what we call a "raw quantum speedup". However, it is essential to realise that the time used by the QPU and measured by the QPT refers only to a subset of the processing required to solve a given problem instance using a quantum annealer. Specifically, a complete quantum algorithm for a problem instance $P$ involves, as a minimum requirement, the following steps:

1. *Conversion:* The problem instance $P$ must be converted into a QUBO instance $Q(P)$, typically via a polynomial-time reduction taking time $t_{\mathrm{conv}}$.

2. *Embedding:* The QUBO problem $Q(P)$ must be embedded into the Chimera hardware graph taking time $t_{\mathrm{embed}}$.

3. *Pre-processing:* The embedded problem is pre-processed, which involves calculating (appropriately scaled) weights for the embedded QUBO problem, taking time $t_{\mathrm{pre}}$.

4. *Quantum processing:* The annealing process is performed on the QPU, taking time $t_{\mathrm{proc}}$.

5. *Post-processing:* The samples are post-processed to choose the best candidate solution, check its validity, and perform any other post-processing methods to improve the solution quality [19, 25] taking time $t_{\mathrm{post}}$. The QUBO solution must finally be converted back to a solution for the original problem $P$.

The total processing time is thus

$$T_Q = t_{\mathrm{conv}} + t_{\mathrm{embed}} + t_{\mathrm{pre}} + t_{\mathrm{proc}} + t_{\mathrm{post}}. \tag{1}$$

The realisation that these other steps must be included in the analysis is emphasised by the fact that in practical problems the embedding time often dominates the time used by the annealer itself. Previous investigations have largely avoided this by focusing on artificial problems "planted" in the Chimera graph so that no embedding is necessary [3, 12, 16, 26, 27]. Although finding a raw speedup in such situations is clearly a necessary condition for a quantum speedup, it is not sufficient for it to be present in practical problems.

To properly benchmark quantum annealing it is necessary to also compare fairly the quantum annealer to a suitable classical algorithm. Indeed, much of the controversy regarding potential speedups with quantum annealing has been due to the fact that quantum annealers have been compared against simulated annealing or simulated quantum annealing. Although such studies certainly have merit and such a speedup is certainly a necessary condition for a real quantum speedup, it has repeatedly been pointed out that classical annealing techniques are generally far from optimal and any observed speedups have disappeared when better classical algorithms were used [12]. In [26], this type of quantum speedup has been termed a "limited speedup". Ideally, one should instead compare annealing against the *best available* classical algorithm for the problem to find a "potential quantum speedup".

# 3 Hybrid quantum-classical computing

Much of the previous effort towards determining whether or not quantum annealing can, in practice, provide a computational speedup has focused on determining the existence of a *raw quantum speedup*, which does not take into account the associated classical processing that is inseparable from a quantum annealer. Such a raw speedup is certainly a necessary condition for practical quantum computational gains, and its study is therefore well justified. However, even if there is a raw speedup there are many reasons why this might not translate into a *practical* quantum speedup.

A practical speedup is possible for a problem if we are able to give a quantum algorithm such that $T_Q < T_C$, where $T_C$ is the classical processing time for the best available classical algorithm for the problem. From the definition of $T_Q$ in (1), it is clear than, even if $t_{\text{proc}} < T_C$, the conversion, embedding and pre/post-processing may provide obstructions to obtaining a practical speedup. In practical terms, the pre- and post-processing tend to add relatively minor (or controllable) overheads, but the conversion and embedding costs pose more fundamental problems.

These difficulties in turning a raw quantum speedup into a practical advantage for practical problems have led to significant interest in "hybrid classical-quantum" approaches (also called "quassical" computations by Allen [4]). Hopefully, combining quantum annealing with classical algorithms may allow otherwise inaccessible speedups to be exploited. Several such hybrid approaches have aimed to overcome the resource limitation arising from the fact that practical problems typically require more qubits than are available on existing devices (as a result of the expansion in number of variables during the conversion stage discussed above) [24, 29].

## 3.1 Hybrid computing to negate embedding costs

Although hybrid approaches have also looked at improving the robustness and quality of embeddings [31], to the best of our knowledge such approaches have not been used to try and mitigate the cost of performing the embedding itself, which, we recall, is often prohibitive to any speedup. In this paper we propose a general hybrid approach to tackle precisely this problem. In particular we aim to show how a raw speedup that is negated by the embedding time (i.e., in particular when $t_{\text{proc}} < T_C$ but $T_Q > T_C$) can nonetheless be exploited to give a practical speedup to certain computational problems.

Our approach is motivated by another hybrid quantum-classical algorithmic proposal which predates the rise of quantum annealing and was introduced with the aim of exploiting Grover's algorithm—the well-known black-box algorithm for quantum unordered database search [15]—in practical applications [22]. The crucial condition for such a problem to be amenable to this hybrid approach is that *the repeated calls to the quantum annealer should be made with the same logical graph embedding*, or at least *permit an efficient method to construct the embedding for one call from the previous ones.* If this condition is satisfied, the cost of the embedding, $t_{\text{embed}}$, can thus be spread out over the several calls, allowing a raw quantum speedup to be exploited.

In order to see how this hybrid approach can help exploit a quantum speedup, we will consider the following general description of a quantum annealing algorithm based on the hybrid approach described above (a more precise analysis would necessarily depend in part on the algorithm in question): some initial classical processing is performed, the embedding of a logical graph into the physical graph is

computed, $m$ instances of a QUBO problem are solved on a quantum annealer, with some classical pre- and post-processing occurring between instances, and some final classical computation is optionally performed. More formally, let us call the overall problem the hybrid algorithm solves $R$, and the $m$ problem instances that must be solved to do so, $P_1, \ldots, P_m$. Recall that the time to solve a single instance $P_i$ on an annealer is $T_Q(P_i)$. As we noted earlier this is, in practical situations, generally dominated by the cost of the embedding and the quantum processing, so $T_Q(P_i)$ can be approximated, for simplicity, as

$$T_Q(P_i) = t_{\text{conv}}(P_i) + t_{\text{embed}}(P_i) + t_{\text{pre}}(P_i) + t_{\text{proc}}(P_i) + t_{\text{post}}(P_i) \approx t_{\text{embed}}(P_i) + t_{\text{proc}}(P_i), \qquad (2)$$

where we have explicitly included the dependence on the problem instance. The hybrid algorithm will thus take time

$$T_H(R) \approx t_1(R) + t_{\text{embed}}(P_1) + \sum_i \left( t_{\text{proc}}(P_i) + t_2(P_i) \right) \approx t_1(R) + t_{\text{embed}}(P_1) + \sum_i t_{\text{proc}}(P_i),$$

where $t_1(R)$ encapsulates any initial and final classical processing associated with combining the solutions $P_i$, and $t_2(P_i)$ is the time of the classical calculation associated with each iteration, which we have assumed to be small compared to $t_{\text{proc}}(P_i)$ since this should simply encompass minor pre- and post-processing between annealing runs, and thus be negligible if the problem is amenable to the hybrid approach. Note that we have made use of the assumption that $t_{\text{embed}}(P_1) \approx t_{\text{embed}}(P_i)$ for $i > 1$, which is a criterion in the suitability of a problem for this hybrid approach.

We note immediately that a standard approach with a quantum annealer, performing the embedding for each instance $P_i$, would take time $T_{\text{std}}(R) = t_1(R) + \sum_i \left( t_{\text{embed}}(P_i) + t_{\text{proc}}(P_i) \right)$. Thus, since in practice $t_{\text{embed}}$ is comparable, if not larger, than $t_{\text{proc}}$, we already have $T_H(R) \ll T_{\text{std}}(R)$. Although this conclusion may seem somewhat trivial, it is important in that it shows already how annealing can provide much larger practical gains for such complex algorithmic problems. More importantly, *it may allow a raw quantum speedup to be exploited practically.*

It is, of course, possible that for certain problems a much more efficient classical algorithm exists for solving $R$ when $m$ is large enough (e.g., there might be an efficient way to map solutions of $P_i$ to $P_j$). Such problems are thus not suitable for such a hybrid approach, and so are not of particular interest to us. Nonetheless, generally a classical algorithm for $R$ may be more intelligent than the standard approach as certain, presumably minor, parts of the computation are likely to be common to solving several $P_i$. Specifically, we can thus rewrite $T_C(P_i) = t_3(P_i) + t_4(P_i)$, where $t_3$ is small compared to $t_4$. The best classical algorithm can then, rather generally, be considered to take time

$$T_C^{\text{best}}(R) = t_5(R) + t_3(P_1) + \sum_i t_4(P_i) = t_6(R) + \sum_i t_4(P_i),$$

where $t_6(R) = t_5(R) + t_3(P_1)$. Crucially, unless the raw quantum speedup is small, we will also have $t_{\text{proc}}(P_i) < t_4(P_i)$. It is thus easy to see that, *for large enough $m$ (i.e., number of $P_i$ to be solved), we will have $T_H(R) < T_C^{\text{best}}(R)$,* and thus the raw quantum speedup will translate into an absolute speedup for the hybrid algorithm.

## 4   Case study: Dynamically weighted maximum-weight independent set

To illustrate the proposed hybrid approach, we discuss in detail a concrete example both from a theoretical and experimental viewpoint.

## 4.1   (Dynamically weighted) Maximum-weight independent set

The core of the problem is the maximum-weighted independent (MWIS) set problem. Recall that an *independent set* $V'$ of vertices of a graph $G = (V, E)$ is a set $V' \subseteq V$ such that for all $\{u, v\} \in E$ we have $\{u, v\} \not\subseteq V'$.

**Maximum-Weight Independent Set (MWIS) Problem**:

*Input:*   A graph $G = (V, E)$ with positive vertex weights $w : V \to \mathbb{R}^+$.
*Task:*   Find an independent set $V' \subseteq V$ such that maximises $\sum_{v \in V'} w(v)$
over all independent sets of $G$.

The general MWIS problem is NP-hard since it encompasses, by restriction, the well-studied non-weighted version [14]. One should note, however, that for graphs of bounded tree-width, the MWIS problem is polynomial-time solvable using standard dynamic programming techniques (see [23]).

Although the MWIS can be readily transformed into a QUBO problem (as we show below), by itself it is not directly suitable for the hybrid approach we proposed. However, a simple variation that we propose here is indeed suitable.

**Dynamically Weighted Maximum-Weight Independent Set (DWMWIS) Problem**:

*Input:*   A graph $G = (V, E)$ with a set of weight functions $W = \{w_1, w_2, \ldots, w_m\}$
where $w_i : V \to \mathbb{R}^+$ for $1 \le i \le m$.
*Task:*   Find independent sets $V_i \subseteq V$ that maximise $\sum_{v \in V_i} w_i(v)$ for each $1 \le i \le m$.

This problem is to solve the MWIS problem on $G$ for each of the $m$ weight assignments $w_i \in W$. For $m = 1$ we obtain again the MWIS problem, but for larger $m$ the problem is suitable for our hybrid approach.

## 4.2   Quantum solution

We now provide a QUBO formulation for the MWIS Problem. Fix an input graph $G = (V, E)$ with positive vertex weights $w : V \to \mathbb{R}^+$. Let $W = \max\{w(v) \mid v \in V\}$ and let $S > W$ be a "penalty weight". We build a QUBO matrix of dimension $n = |V|$ such that:

$$
Q_{(i,j)} = \begin{cases} 0, & \text{if } i > j \text{ or } \{i, j\} \notin E, \\ -w(v_i), & \text{if } i = j, \\ S, & \text{if } i < j \text{ and } \{i, j\} \in E. \end{cases} \tag{3}
$$

**Theorem 1** ([2]). *The QUBO formulation given in* (3) *solves the MWIS Problem.*

In order to adapt the MWIS solution above to the DWMWIS problem, note that the non-zero entries of the QUBO formulation (3) depend only on the structure of the graph and not on the weight function $w$. Thus, in order to solve the DWMWIS problem, for each weight assignment $w_i$ the same embedding of the graph into the D-Wave physical graph can be used, meaning that a hybrid algorithm based around the MWIS solution above can readily be implemented. More specifically, following the hybrid algorithm described in Section 3.1 for instances $P_1, \ldots, P_m$ (where each $P_i$ uses weight function $w_i$), we perform

the embedding once (entailing a time $t_{\text{embed}}(P_1)$) and then solve the MWIS problem for each weight assignment $w_i$ (taking times $t_{\text{proc}}(P_i)$) using the QUBO solution outlined above.

## 4.3   Classical baseline

The main objective of studying the DWMWIS example in detail is to exhibit experimentally the advantage that the hybrid approach can provide over a standard annealing-based approach. Nonetheless, it is helpful to further compare this to the performance of a classical baseline algorithm for comparison and to help highlight this advantage, even if we do not necessarily expect to see an absolute quantum speedup from the hybrid algorithm.

To this end, for a given input graph $G = (V, E)$ with positive vertex weights $w : V \to \mathbb{R}^+$, we construct a Binary Integer Programming (BIP) instance with $n = |V|$ binary variables as follows. To each vertex $v_i$ in $G$ we associate the binary variable $x_i$, and for notational simplicity we will denote the collection of variables $x_i$ by a binary vector $\mathbf{x} = (x_0, x_1, \cdots, x_{n-1})$. We thus have the BIP problem instance:

$$
\begin{aligned}
&\text{maximise} \ \sum_{v_i \in V} w(v_i) x_i \\
&\text{subject to } x_i + x_j \leq 1 \text{ for all } \{v_i, v_j\} \in E.
\end{aligned}
\tag{4}
$$

Each constraint in (4) enforces the property that no adjacent vertices are chosen in the independent set while the objective function ensures an independent set with maximum sum value is chosen. Assuming we have the binary vector $\mathbf{x}$ which yields the optimal value of objective function (4), we take $D(\mathbf{x}) = \{v_i \mid x_i = 1\}$ to be the set of vertices selected as the maximum weighted independent set.

**Theorem 2** ([2]). *The BIP formulation given in (4) solves the MWIS problem.*

The classical baseline used in the analysis is based on an implementation of the BIP formulation in Sage Math [28], which has a well developed and optimised Mixed Integer Programming library. To ensure that a fair comparison with the hybrid algorithm is possible, we formulate the classical algorithm for the overall DWMWIS problem such that *the set of constraints in the BIP formulation is only computed once*.

## 4.4   Experimental definition and procedure

To study the performance of the hybrid DWMIWS algorithm in a practical setting, we made use of a D-Wave 2X quantum annealer with 1098 active physical qubits [8] to compare the performance of three algorithms on a selection DWMWIS problem instances: the "standard" quantum algorithm, in which the embedding is re-performed for each weight assignment; the hybrid DWMWIS algorithm; and the classical BIP-based solution described above. We present here a summary of the experimental procedure and results; a more detailed presentation and analysis is available in an extended version of the paper [2].

To this end we analyse the algorithms on a range of different graphs, in particular choosing 155 graphs from a variety of common graph families with between 2 and 126 vertices. Each graph was used to generate a single DWMWIS problem instance with $m = 100$ weight assignments, each randomly generated as floating point numbers rounded to 2 decimal places within the range $[0.0, 1.0]$. The problem instances were generated as standard adjacency list representations using SageMath [28] with random weights.

The same procedure is used for the "standard" quantum algorithm, except the cost of the embedding is incurred for each weight assignment.

Since we are primarily interested in negating the impact of the embedding process in general applications, we made use of D-Wave's heuristic embedding algorithm [10] to embed each logical graph in the physical graph. Each graph was embedded 10 times to estimate $t_{\text{embed}}$ for each problem instance. Finally, our tests were run with D-Wave's post-processing optimisation enabled. While this adds a small overhead in time, this is well within the spirit of hybrid quantum-classical computing, and allowed us to solve more problems. This post-processing method processes small batches of samples while the next batch is being processed [9]. This ensures that it only contributes a constant overhead in time for each MWIS problem instance *independent of the number of samples (and thus $k_{99}$).*

## 4.5   Results and analysis

For each DWMWIS problem instance (i.e., for each graph $G$) the times $T_H$ and $T_{\text{std}}$ were calculated, following the approach described in Section 3.1, as

$$T_H = t_{\text{embed}} + \sum_i \left( t_{\text{prog}}(P_i) + k_{99}(P_i)t_{\text{anneal}} + t_{\text{post}}(P_i) \right),$$

$$T_{\text{std}} = \sum_i \left( t_{\text{embed}} + t_{\text{prog}}(P_i) + k_{99}(P_i)t_{\text{anneal}} + t_{\text{post}}(P_i) \right),$$

where $k_{99}(P_i)$ is the $k_{99}$ value for weight assignment $w_i$ and $t_{\text{anneal}} = 309\mu s$. Both $t_{\text{prog}}(P_i)$ and $t_{\text{post}}(P_i)$ are of the order of 20ms. Note that the processing time $t_{\text{proc}}$ defined earlier is, for this approach to the DWMWIS problem, given by $t_{\text{proc}} = t_{\text{prog}}(P_i) + k_{99}(P_i)t_{\text{anneal}} + t_{\text{post}}(P_i)$. The classical time $T_C$ was taken as the processor time for the classical algorithm described above. The results are summarised in Figures 1(a) and 1(b), which show how the hybrid times $T_H$ compare to both $T_{\text{std}}$ and $T_C$. Error bars are calculated from the observed variation in $t_{\text{embed}}$, the number of optimal solutions found $N_{\text{opt}}$, and the post-processing time $t_{\text{post}}$. Of these, the error in $t_{\text{post}}$ is the dominant factor, and largely arises from the uncontrollability of the post-processing environment, which is performed remotely within the D-Wave processing pipeline. However, this variation did not result in any significant variation in success probability of the annealing, so it seems the amount of post-processing performed was constant.

First and foremost, from the results shown in Figure 1(a) the extent of the advantage of the hybrid approach is evident. Indeed, this is to be expected given that, for a given DWMWIS problem, they differ (by definition) by $99 \times t_{\text{embed}}$. Although this might seem a trivial confirmation of this fact, the results help illustrate the extent of the advantage that the hybrid approach can have for such problems, a consequence of the absolute cost of the embedding. This is visible in Figure 2, showing $t_{\text{embed}}$ as a function of the number of vertices in a graph.

From Figure 1(b) it is also evident that no absolute quantum speedup was observed using the hybrid algorithm, and indeed there is a vast difference in scale between $T_C$ and $T_H$: the "hardest" problem was solved classically in less than 200ms, whereas the hybrid algorithm required almost 60 times as much time to solve it correctly. The inability to observe any raw speedup is hardly surprising when one notes that, even if $k_{99} = 1$ and $t_{\text{embed}} = t_{\text{post}} = 0$, the fact that $t_{\text{prog}} \approx 20$ms means that that one would have $T_H > 2000$ms. This programming time thus adds an essentially constant overhead, which would have less of an impact as larger problems (for which $k_{99}$ is much larger) become solvable.
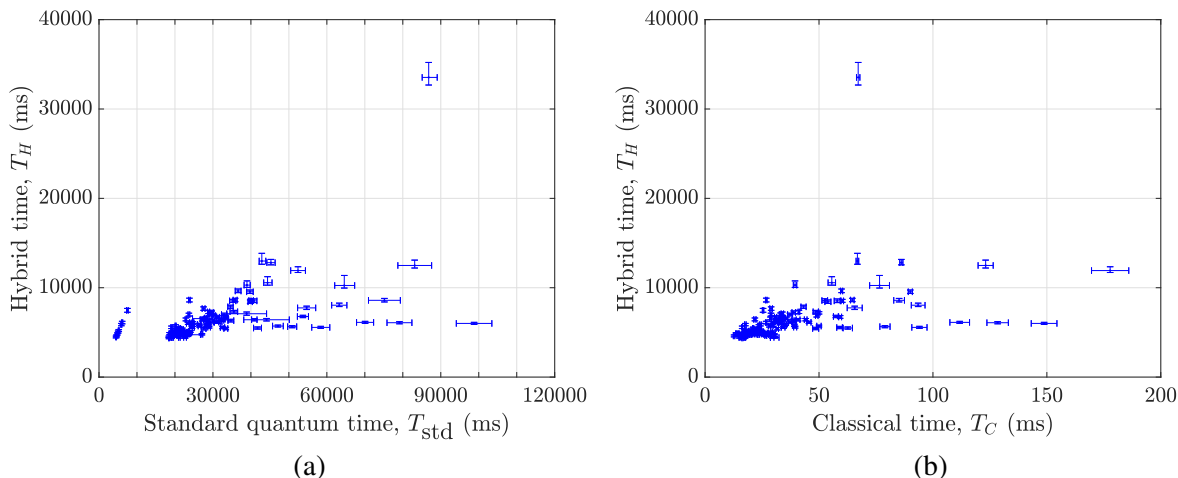
Figure 1: Plots of (a) an upper bound for $T_{\text{std}}$ against $T_H$; and (b) $T_C$ against $T_H$ for each DWMWIS problem instance. All times are in ms.

Despite the absence of no overall speedup, it is interesting to examine the scaling behaviour of the hybrid approach, for which it will be useful to consider the "classical speedup ratio" $R_C = T_H/T_C$. In Figure 3 we show the scaling behaviour of $R_C$ against two reasonable proxies of problem difficulty: the graph order $|V|$, which is proportional to the problem size, and the classical time $T_C$. While there is much uncertainty in the exact nature of the scaling, this results indicate that the Hybrid algorithm has a better scaling behaviour than the classical algorithm. This is more evident in Figure 4, where $R_C$ is plotted for specific graph families. Thus, although no quantum speedup was found, the results leave open the possibility that such a speedup will be attainable in the future on larger devices with better control of the qubits, although many unknowns may plausibly alter the scaling behaviour in the future.

Nevertheless, the experiment was a successful proof-of-concept for the hybrid paradigm we have presented. In particular, the hybrid algorithm we implemented provided large absolute gains over the standard quantum approach and showed good scaling behaviour. As larger and more efficient devices become available and more problems of practical interest are studied, it will become clearer if/when a quantum speedup might be obtainable in practise.

## 5   Conclusion

In this paper we presented a hybrid quantum-classical paradigm for quantum annealing algorithms aimed at countering the significant cost of the embedding process. This approach is not only a hybrid paradigm but serves equally as a guide to identifying problems that may be amenable to quantum annealing. In particular, we identify those problems that require solving a large number of related subproblems, each of which can be directed solved via annealing, may permit a hybrid approach. This is obtained by reusing and modifying embeddings for the related subproblems.
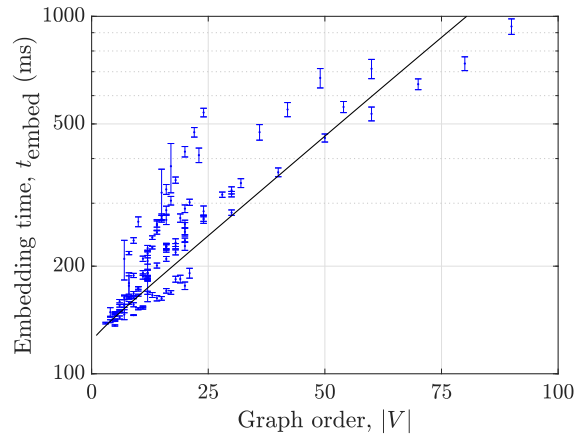
Figure 2: Plot of graph order $|V|$ against the embedding time $t_{\text{embed}}$. Note the logarithmic scale in time.
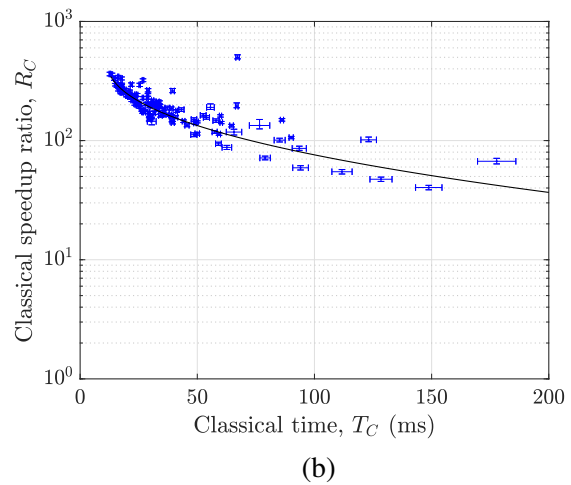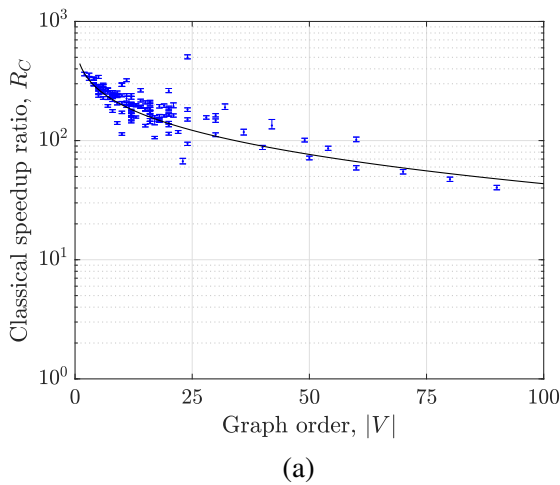


Figure 3: Logarithmic plots of the scaling behaviour of the classical speedup ratio $R_C$ for the DWMWIS problem instances: (a) graph order $|V|$ against $R_C$; and (b) classical time $T_C$ against $R_C$.
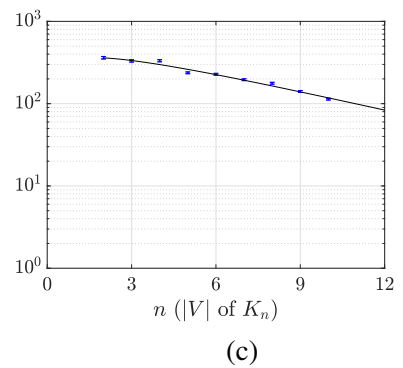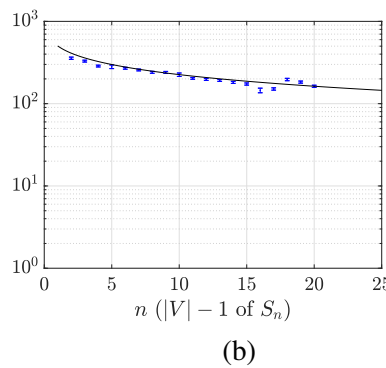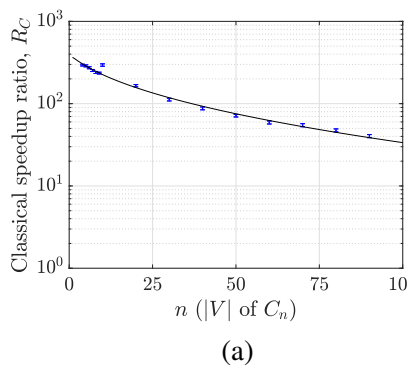


Figure 4: Plots of the classical speedup ratio $R_C$ against $n$ for three families of graphs parameterised by $n$: (a) the $C_n$ graphs; (b) the $S_n$ graphs; (c) the $K_n$ graphs.

Our hybrid approach, along with its successful proof-of-principle implementation, sets the groundwork for addressing more complex problems of practical interest. Choosing correctly suitable problems is a major step in finding practical uses for quantum computers in the near term future, and with deft choices, quantum speedups from hybrid approaches might soon be realisable.

# References

[1] S. Aaronson (2010): *BQP and the polynomial hierarchy*. In: *STOC '10 Proceedings of the forty-second ACM symposium on Theory of computing*, p. 141, doi:10.1145/1806689.1806711.

[2] A. A. Abbott, C. S. Calude, M. J. Dinneen & R. Hua (2018): *A Hybrid Quantum-Classical Paradigm to Mitigate Embedding Costs in Quantum Annealing*. CDMTCS Research Report Series 520.

[3] S. Boixo, T. F. Rønnow, S. V. Isakov, Z. Wang, D. Wecker, D. A. Lidar, J. M. Martinis & M. Troyer (2014): *Evidence for quantum annealing with more than one hundred qubits*. Nat. Phys. 10, p. 218, doi:10.1038/nphys2900.

[4] C. S. Calude, E. Calude & M. J. Dinneen (2015): *Adiabatic Quantum Computing Challenges*. ACM SIGACT News 46(1), p. 40, doi:10.1145/2744447.2744459.

[5] C. S. Calude & M. J. Dinneen (2016): *Solving the Broadcast Time Problem Using a D-Wave Quantum Computer*. In A. Adamatzky, editor: *Advances in Unconventional Computing*, chapter 17, *Emergence, Complexity and Computation* 22, Springer International, Switzerland, p. 439, doi:10.1007/978-3-319-33924-5_17.

[6] A. Cho (2014): *Quantum or not, controversial computer yields no speedup*. Science 344, p. 1330, doi:10.1126/science.344.6190.1330.

[7] V. Choi (2008): *Minor-embedding in adiabatic quantum computation: I. The parameter setting problem*. Quantum Inf. Processing 7, p. 193, doi:10.1007/s11128-008-0082-9.

[8] D-Wave Systems Inc. (2016): *The D-Wave 2X™ Quantum Computer Technology Overview*. Available at `http://www.dwavesys.com/sites/default/files/D-Wave%202X%20Tech%20Collateral_1016F_0.pdf`.

[9] D-Wave Systems Inc. (2016): *Postprocessing Methods on D-Wave Systems*. Tech. Report Release 2.4 09-1105A-B.

[10] D-Wave Systems Inc. (2017): *Programming with QUBOs*. Tech. Report Release 2.4 09-1002A-C.

[11] D-Wave Systems Inc. (2017): *The D-Wave 2000Q™ Quantum Computer Technology Overview*. Available at `https://www.dwavesys.com/sites/default/files/D-Wave%202000Q%20Tech%20Collateral_0117F2.pdf`.

[12] V. S. Denchev, S. Boixo, S. V. Isakov, N. Ding, R. Babbush, V. Smelyanskiy, J. Martinis, & H. Neven (2016): *What is the Computational Value of Finite-Range Tunneling?* Phys. Rev. X 6, p. 031015, doi:10.1103/PhysRevX.6.031015.

[13] E. Farhi, J. Goldstone, S. Gutman & M. Sipser (2000): *Quantum Computation by Adiabatic Evolution*. arXiv:quant-ph/0001106.

[14] M. R. Garey & D. S. Johnson (1979): *Computers and Intractability. A Guide to the Theory of NP-Completeness*. Freeman, San Francisco.

[15] L. K. Grover (1996): *A fast quantum mechanical algorithm for database search*. In: *Proceedings, 28th Annual ACM Symposium on the Theory of Computing (STOC)*, p. 212, doi:10.1145/237814.237866.

[16] I. Hen, J. Job, T. Albash, T. F. Rønnow, M. Troyer & D. A. Lidar (2015): *Probing for quantum speedup in spin-glass problems with planted solutions*. *Phys. Rev. A* 92, p. 042325, doi:10.1103/PhysRevA.92.042325.

[17] M. W. Johnson, M. H. S. Amin, S. Gildert, T. Lanting, F. Hamze, N. Dickson, R. Harris, A. J. Berkley, J. Johansson, P. Bunyk, E. M. Chapple, C. Enderud, J. P. Hilton, K. Karimi, E. Ladizinsky, N. Ladizinsky, T. Oh, I. Perminov, C. Rich, M. C. Thom, E. Tolkacheva, C. J. S. Truncik, S. Uchaikin, J. Wang, B. Wilson & G. Rose (2011): *Quantum annealing with manufactured spins*. *Nature* 473, p. 194, doi:10.1038/nature10012.

[18] A. D. King, T. Lanting & R. Harris (2015): *Performance of a quantum annealer on range-limited constraint satisfaction problems*. arXiv:1502.02098 [quant-ph].

[19] A. D. King & C. C. McGeoch (2014): *Algorithm engineering for a quantum annealing platform*. arXiv:1410.2628 [cs.DS].

[20] J. King, S. Yarkoni, M. M. Nevisi, J. P. Hilton & C. C. McGeoch (2015): *Benchmarking a quantum annealing processor with the time-to-target metric*. arXiv:1508.05087 [quant-ph].

[21] T. D. Ladd, F. Jelezko, R. Laflamme, Y. Nakamura, C. Monroe & J. L. O'Brien (2010): *Quantum computers*. *Nature* 464, p. 45, doi:10.1038/nature08812.

[22] M. Lanzagorta & J. K. Uhlmann (2005): *Hybrid quantum-classical computing with applications to computer graphics*. In: *ACM SIGGRAPH 2005 Courses*, SIGGRAPH '05, ACM, New York, NY, doi:10.1145/1198555.1198723.

[23] D. Marx (2010): *Fixed parameter algorithms. Part 2: Treewidth*. Available at `http://www.cs.bme.hu/~dmarx/papers/marx-warsaw-fpt2`. Open lectures for PhD students in computer science, University of Warsaw, Poland.

[24] J. R. McClean, J. Romero, R. Babbush & A. Aspuru-Guzik (2016): *The theory of variational hybrid quantum-classical algorithms*. *New J. Phys.* 18, p. 023023, doi:10.1088/1367-2630/18/2/023023.

[25] K. L. Pudenz (2016): *Parameter Setting for Quantum Annealers*. In: *20th IEEE High Performance Embedded Computing Workshop Proceedings*, doi:10.1109/HPEC.2016.7761619.

[26] T. F. Rønnow, Z. Wang, J. Job, S. Boixo, S. V. Isakov, D. Wecker, J. M. Martinis, D. A. Lidar & M. Troyer (2014): *Defining and detecting quantum speedup*. *Science* 345, p. 420, doi:10.1126/science.1252319.

[27] S. W. Shin, G. Smith, J. A. Smolin & U. Vazirani (2014): *How "Quantum" is the D-Wave Machine?* arXiv:1401.7087 [quant-ph].

[28] The Sage Developers (2017): *SageMath, the Sage Mathematics Software System (Version 8.0)*. Available at `http://www.sagemath.org`.

[29] T. T. Tran, M. Do, E. G. Rieffel, J. Frank, Z. Wang, B. O'Gorman, D. Venturelli & J. C. Beck (2016): *A Hybrid Quantum-Classical Approach to Solving Scheduling Problems*. In: *Proceedings of the Ninth International Symposium on Combinatorial Search*, AAAI.

[30] D. Venturelli, D. J. J. Marchand & G. Rojo (2015): *Job Shop Scheduling Solver based on Quantum Annealing*. arXiv:1506.08479 [quant-ph].

[31] W. Vinci, T. Albash, G. Paz-Silva, I. Hen & D. A. Lidar (2015): *Quantum annealing correction with minor embedding*. *Phys. Rev. A* 92, p. 042310, doi:10.1103/PhysRevA.92.042310.

[32] S. Yarkoni, A. Plaat & T. Bäck (2017): *First results solving arbitrarily structured Maximum Independent Set problems using quantum annealing*. Available at `http://liacs.leidenuniv.nl/~plaata1/papers/MIS_yarkoni.pdf`.