

This is the accepted manuscript made available via CHORUS. The article has been published as:

Thermalization, Freeze-out, and Noise: Deciphering Experimental Quantum Annealers

Jeffrey Marshall, Eleanor G. Rieffel, and Itay Hen

Phys. Rev. Applied **8**, 064025 — Published 26 December 2017

DOI: [10.1103/PhysRevApplied.8.064025](https://doi.org/10.1103/PhysRevApplied.8.064025)

Thermalization, Freeze-out and Noise: Deciphering Experimental Quantum Annealers

Jeffrey Marshall,¹ Eleanor G. Rieffel,² and Itay Hen^{1,3}

¹*Department of Physics and Astronomy, and Center for Quantum Information Science & Technology,
University of Southern California, Los Angeles, California 90089, USA*

²*QuAIL, NASA Ames Research Center, Moffett Field, California 94035, USA*

³*Information Sciences Institute, University of Southern California, Marina del Rey, California 90292, USA*

By contrasting the performance of two quantum annealers operating at different temperatures, we address recent questions related to the role of temperature in these devices and their function as ‘Boltzmann samplers’. Using a method to reliably calculate the degeneracies of the energy levels of large-scale spin-glass instances, we are able to estimate the instance-dependent effective temperature from the output of annealing runs. Our results corroborate the ‘freeze-out’ picture which posits two regimes, one in which the final state corresponds to a Boltzmann distribution of the final Hamiltonian with a well-defined ‘effective temperature’ determined at a freeze-out point late in the anneal, and another regime in which such a distribution is not necessarily expected. We find that the output distributions of the annealers do not in general correspond to a classical Boltzmann distribution for the final Hamiltonian. We also find that the effective temperatures at different programming cycles fluctuate greatly, with the effect worsening with problem size. We discuss the implications of our results for the design of future quantum annealers to act as more effective Boltzmann samplers and for the programming of such annealers.

I. INTRODUCTION

A handful of recent studies suggest that quantum annealers may be well suited to function as fast thermal samplers [1–4]. By taking advantage of their finite temperature nature [3–7], potentially they may sample from Boltzmann distributions of certain cost functions more efficiently than can be done classically. Such a capability opens up the exciting possibility of applications of quantum annealing to so-far-uncharted avenues of research, with immediate applications to domains such as deep learning networks and restricted Boltzmann machines [2, 3, 8].

The main mechanisms that determine the distributions from which output configurations are drawn are thus far unclear. Further insights into the role of temperature, and the capabilities of experimental quantum annealing optimizers to quickly thermalize, are challenging to obtain due to the limited ability to probe the inner workings of these machines, as well as the lack of control over most operating parameters [3, 4, 8].

To circumvent these difficulties, we devised an experiment, directly comparing the performance of *two* commercially available quantum annealers operating at different temperatures (we shall refer to those as ‘hot’ and ‘cold’ henceforth). This key difference, together with a newly devised method to accurately calculate the degeneracies of certain large-scale spin-glass instances, offers us a unique opportunity to study the effects of temperature. Our results indicate that these instances do not in general equilibrate at Boltzmann distributions corresponding to the final classical Hamiltonian, but are significantly affected by nonzero quantum fluctuations and noise. Our results corroborate the ‘freeze-out’ picture [1, 2, 9], which posits one regime in which the final state corresponds to a Boltzmann distribution of the final Hamiltonian with well-defined ‘effective (classical) temperature’ de-

termined at a generally unknown freeze-out point late in the anneal, and another regime in which such a distribution would not necessarily be expected. While providing evidence for this picture, our results speak against the hypothesis that most instances fall in the first regime.

We find that these effective temperatures fluctuate greatly at different programming cycles, with the effect worsening with problem size. We discuss factors potentially contributing to this adverse effect, including so-called J -chaos in which control errors and other sources of noise mean that the problem run on the machine is different from the one programmed in. We discuss the implications of our results for the design of future quantum annealers to act as efficient Boltzmann samplers and for the programming of such annealers.

A. Quantum annealing and quantum annealers

Standard transverse field quantum annealing works by evolving the system over rescaled time $s = t/\mathcal{T} \in [0, 1]$ where t is time and \mathcal{T} is the overall runtime of the annealing process. The total Hamiltonian of the system is given by

$$H(s) = A(s)H_d + B(s)H_p, \quad (1)$$

where $H_p = \sum_{\langle i,j \rangle} J_{ij} \sigma_i^z \sigma_j^z + \sum_i h_i \sigma_i^z$ is the programmable Ising spin-glass problem (the final Hamiltonian) to be sampled defined by the parameters $\{J_{ij}, h_i\}$, and $H_d = -\sum_i \sigma_i^x$ is a transverse-field Hamiltonian which provides the quantum fluctuations (the initial Hamiltonian). We identify two dimensionless scales associated with the annealing, namely, the one associated with quantum fluctuations $Q(s) = A(s)/B(s)$ and the scale associated with thermal fluctuations $k_B T/B(s)$,

both of which are shown in Fig. 1 for both the ‘hot’ and ‘cold’ processors.

Current quantum annealers suffer from intrinsic control errors (ICE) [6, 10] such as imperfect digital-to-analog conversion when programming the problem parameters onto the machine, and $1/f$ -noise whose effect is parameter changes even within a single programming cycle (a consecutive batch of anneals run on the machine) [11, 12]. For both contrasted quantum annealers, these random errors may be approximated as normally distributed according to $\sim \mathcal{N}(0, 0.05J)$ [resp. $\sim \mathcal{N}(0, 0.03h)$] where J (resp. h) is the maximal value over all the programmed J_{ij} (resp. h_i). Some problems have resilience to such errors [5, 13], whereas others are susceptible to a phenomenon referred to as J -chaos, in which output ‘solutions’ correspond to the wrong, or malformed, problem, generally reducing the success probability [6, 13–17].

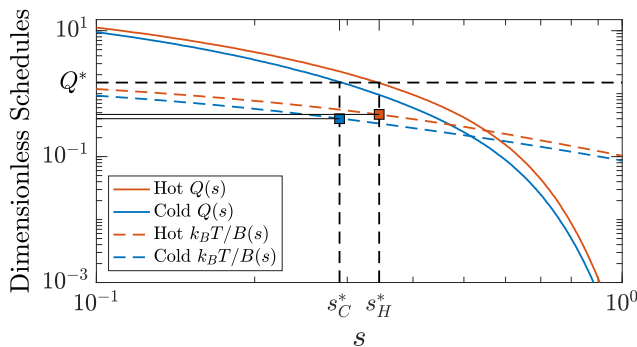


FIG. 1. **Dimensionless annealing schedules and temperatures of the hot and cold DW2 processors.** We plot quantum fluctuations $Q(s)$, and thermal fluctuations $k_B T / B(s)$, as a function of rescaled annealing time, s , for each machine. An example of freeze-out is shown; the distribution approximately halts at some fixed, instance-dependent value of $Q^* := Q(s^*)$ (dashed black line) which corresponds to a freeze-out point s^* (denoted s_C^*, s_H^*) and dimensionless effective temperature $k_B T / B(s^*)$ (squares), for each machine. The solid black lines illustrate that freeze-out may correspond to different effective temperatures.

B. Freeze-out conjecture

If problems thermalized instantly, quantum annealers would return configurations sampled from a Boltzmann distribution, in which each configuration c has weight proportional to $e^{-\beta_{\text{ideal}}^{\text{eff}} E_c}$, where E_c is the configuration’s classical cost (under H_p) and $\beta_{\text{ideal}}^{\text{eff}} \equiv B(1)/k_B T$ is an effective dimensionless inverse temperature, with T being the operating temperature of the machine [18]. It is known, however, that effective inverse-temperatures β^{eff} extracted from experimentally sampled distributions are usually lower than $\beta_{\text{ideal}}^{\text{eff}}$, and that the observed inverse-temperatures differ across problem instances [1–3].

The freeze-out conjecture [1, 2, 9] explains these high observed effective temperatures by positing a “small Q regime” in which the evolution is quasi-static, returning a final population that is close to a Boltzmann distribution of H_p with a well-defined effective temperature, and a regime in which the final population would not necessarily be of this form. In the first regime, the final distribution is determined by a ‘freezing’ of the evolution (after which no dynamics occur) at an unknown, instance-dependent, but physical-temperature independent, ‘freeze-out’ point s^* where thermal fluctuations, whose strength is coupled to the quantum fluctuations $Q(s^*)$ driving the system, become negligibly small [19].

As illustrated in Fig. 1, the freeze-out point is conjectured to happen at a temperature-independent (but instance-dependent) value s^* [1]. Only when $Q(s^*)$ at the freeze-out is small is the final distribution expected to be a classical Boltzmann distribution for H_p with (dimensionless) effective temperature $\beta^{\text{eff}} = B(s^*)/k_B T$; otherwise, the resultant distribution will generally not correspond to an equilibration at any given point, but may instead result from different parts of the system equilibrating at different temperatures and times [1].

C. High-level approach

We proceed by taking as a working hypothesis that most instances have a well-defined freeze-out point in the range $A(s) \ll B(s)$. We work through the implications of this hypothesis, and demonstrate empirically that it does not hold for the majority of instances. We do so by estimating a freeze-out point from the data for each instance, and then checking whether or not that point falls in the $A(s) \ll B(s)$ regime. Most instances fail this consistency check. Outside of that regime, the freeze-out conjecture does not predict a well-defined freeze-out point; different parts of the system may freeze at different times, and even if an instance does have a well-defined freeze-out point outside $A(s) \ll B(s)$, the distribution would have a strong quantum component (from H_d), so would be a distribution of quantum states, and not of the form $e^{-\beta^{\text{eff}} H_p}$. These results are consistent with the freeze-out conjecture, but not with the hypothesis that most instances fall within the freeze-out regime that yields a classical Boltzmann distribution.

II. EXPERIMENT AND METHODS

We make use of two 512-qubit D-Wave Two (DW2) quantum annealers [20]. The mean temperatures of the ‘hot’ and ‘cold’ machines were about 16.0 mK and 13.2 mK, respectively (further details are provided in Appendix A).

We designed 1300 random spin-glass instances of the planted-solution type [21] for each of seven different prob-

lem sizes corresponding to $L \times L$ grids of 8-qubit cells of the hardware DW2 Chimera graph with $L = 2 \dots 8$ (see Fig. 2). We generated these instances as per Ref. [21] (the reader is referred to Appendix B for more details). This class of instances is particularly suitable for our purposes for two main reasons: i) the ground state energies of the generated problems are known in advance, and ii) the *exact* degeneracies of the ground and first excited states are computable [4]. These two facts allow us to, with high accuracy and confidence, measure β^{eff} , as will be explained below. We generated instances on the intersection of the two hardware graphs (501 qubits) in order to avoid biases associated with malfunctioning qubits on either machine (as shown in Fig. 2).

To gather our statistics, each instance was run 440,000 times over 22 ‘programming cycles’ on each machine, with anneal times in range $[20-40] \mu\text{s}$. A programming cycle consists of running the same instance sequentially on a single machine up to (as chosen by the user) 20,000 times, from which statistics are returned; from each programming cycle we obtained the ground state success probability (how often the ground state of the problem was found). We use this data to estimate β^{eff} .

To evaluate β^{eff} we employ two independent, complementary, techniques, which together allow us to estimate with high accuracy and confidence the degeneracies of the energy levels of the problem instances. The first, the well-known Wang-Landau (WL) entropic sampler [22], statistically estimates the degeneracy of the energy levels (see Appendix C for technical details). Since the WL algorithm is prone to statistical errors as well as false convergences, we employ in parallel a newly-devised algorithm that uses the feature that planted-solution instances can be written as a sum of local terms [4]. The algorithm computes the degeneracies of the ground and first excited states *exactly*. When the WL estimate is outside $\pm 5\%$ of the exact value for either the ground or first excited state, we discard this instance as we know it has not converged properly. The combination of these two algorithms allows for the faithful estimation of the degeneracies. We show in Fig. 3 an example of a successful implementation of these two algorithms, where the Wang-Landau ground and first excited estimates are within 5% error of the exact values.

The above procedure yielded some 2200 instances in total, for problem sizes up to 282 qubits, for which we were able to accurately calculate the degeneracies. The difficulty in obtaining an accurate measurement, especially for the larger problems, was due mainly to i) the D-Wave machine not being able to solve many of the ‘hard’ problems, ii) there were too many degenerate states for the exact counter to enumerate (exceeded our chosen cut-off value of 10^7 ground states, which become prohibitively expensive to compute), or iii) Wang-Landau estimate deviated too far from exact counter results (generally from under-sampling the low energy states).

Armed with these degeneracies, we estimate the inverse-temperature β^{eff} for each instance by minimizing

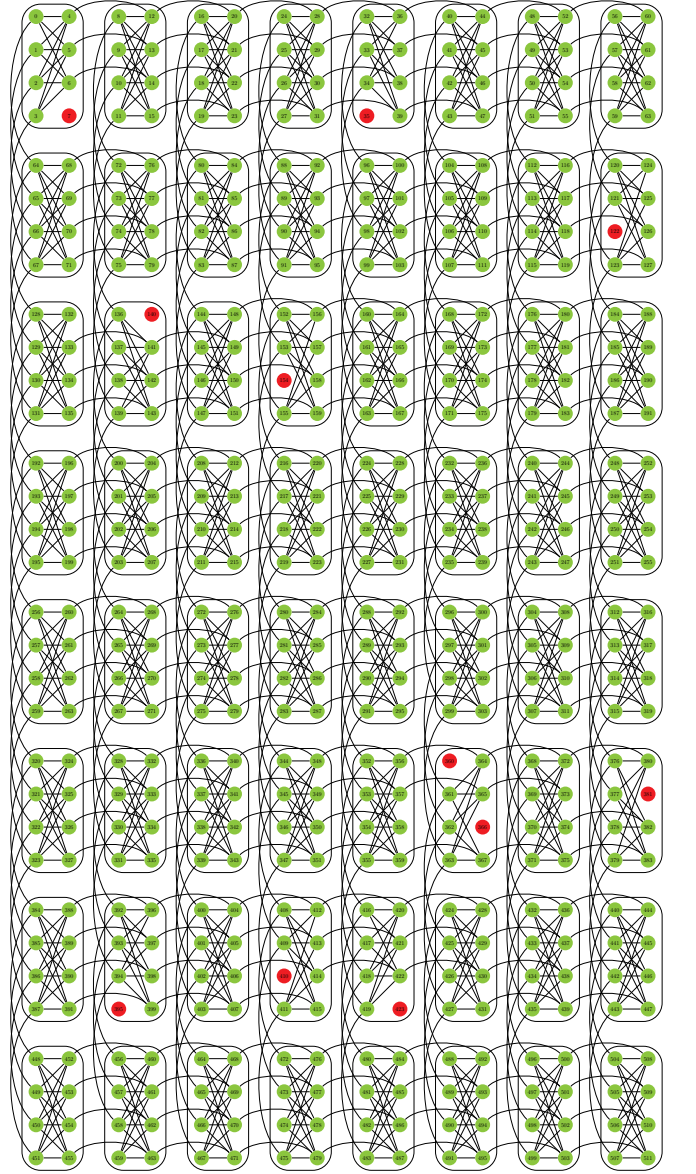


FIG. 2. **Intersected Chimera.** Intersection of the two D-Wave Chimera graphs, with 501 operating qubits. Red disks denote non-operational qubits on one (or both) of the two machines. We denote the size of a square subgraph by an integer $L \leq 8$.

the distance between the observed ground state success probability P_0 and the predicted one (i.e., the conjectured Boltzmann distribution):

$$\left| P_0 - \left(\sum_{k=0} \frac{g_k}{g_0} e^{-\beta^{\text{eff}}(E_k - E_0)} \right)^{-1} \right|. \quad (2)$$

Here, $\{g_k, E_k\}$ are the degeneracy and energy of the k -th level, respectively. The total number of instances for which β^{eff} was successfully estimated, for each problem size $L = 2 \dots 8$, is $[664,745,449,266,38,0,0]$.

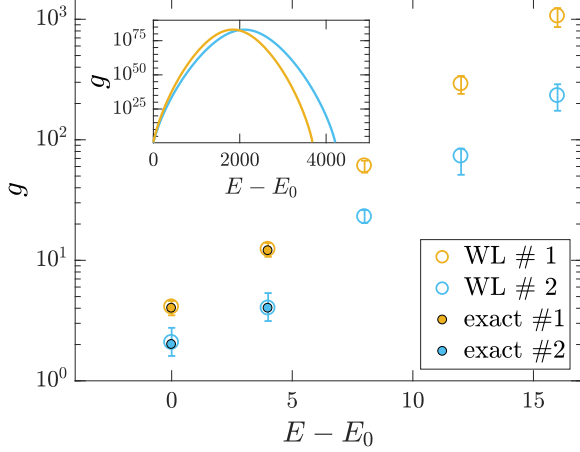


FIG. 3. **Degeneracy counting.** Main figure: The degeneracies of the first five energy levels of two of the 282-qubit problem instances as found by the Wang-Landau algorithm (error bars represent 95% confidence interval). The WL degeneracies of the first two levels lie on top of the computable exact values (solid circles). Inset: Degeneracies of all levels as a function of energy, for the same instances, as obtained by the WL algorithm.

III. RESULTS AND ANALYSIS

A. Thermalization

Figure 4 (top) plots the median inverse temperature β^{eff} for each instance and machine. Error bars indicate the maximum and minimum value of β^{eff} over all programming cycles. Evident is the overall strong linear correlation between the (inverse) effective temperatures of the two machines (Pearson coefficient 0.94). Most instances fall within, or near, the ‘thermal range’ (see caption) predicted by the ratio of physical temperatures of the machines [see yellow band in Fig. 4 (top)], illustrating the key functional role of temperature in the success probability of these problems. If the instances were thermalizing at the end of the anneal, however, we would expect to observe $\beta_{\text{ideal}}^{\text{eff}}$ of 9.7 and 11.7 (shown in Fig. 4) for the hot and cold machines, respectively. Instead, the values we observe are well below this mark: $\beta^{\text{eff}} \in [2, 7]$. Thus, we are finding effective temperatures up to six times higher than would be expected from a simple thermalization picture. Moreover, the median ratio of β^{eff} for the two machines, $R_\beta = 1.11 \pm 0.05$ (95% confidence interval) [23], is well below the ratio of the physical temperatures, $R_\beta^{\text{ideal}} = 1.21 \pm 0.02$ indicating an effective average temperature ratio of about 92% of the ‘thermal’ ratio of $s = 1$. We now examine the extent to which the freeze-out picture can explain these discrepancies.

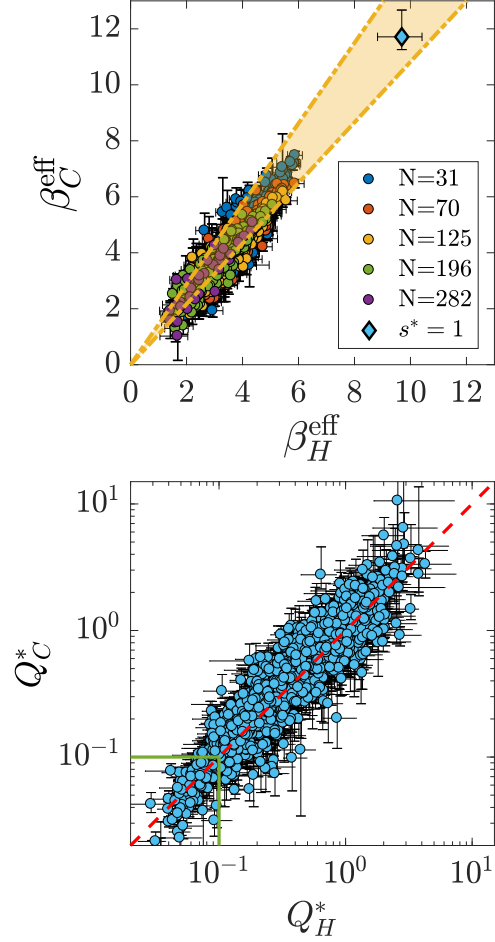


FIG. 4. **Top: Effective inverse temperatures.** Comparison of effective inverse temperatures on the set of instances for both the hot processor, β_H^{eff} , and the cold processor, β_C^{eff} . Error bars represent the highest and lowest values found over all programming cycles. The yellow band represents the range of physical temperature fluctuations between the devices, given by the linear relationship $y = R_\beta^{\text{ideal}} x$ (defined in main text). The blue diamond represents β^{eff} were (classical) thermalization to take place at the end of the anneal. Different problem sizes have different colors (see legend). **Bottom: Quantum fluctuations at freeze-out.** Scatter plot of the extracted $Q^* := Q(s^*)$ strengths (blue) for each instance. Ideally, both machines would yield the same value for each instance (red $y = x$ line). Median ratio is $R_Q^{\text{small}} = 1.01$ for small $Q(s^*)$ values (within the green square). The outputs are strongly correlated (Pearson coefficient 0.92). Error bars represent the range of Q^* over all programming cycles.

B. Freeze-out

While the freeze-out point for each instance is unknown, its temperature independence means the estimates for the freeze out point should be the same whether from the cold machine or hot machine data. Using the estimated $\beta^{\text{eff}} = B(s^*)/k_B T$, from which we can obtain a freeze-out point s^* given the known operating temperatures and annealing schedules, we directly calculate $Q(s^*)$. We then check whether the freeze-out point for each instance is the same for the two machines. We plot $Q(s^*)$ for each instance in Fig. 4 (bottom). For instances with small Q (we take $Q < 10^{-1}$), we find excellent correspondence between the two machines, with an average ratio of $R_Q^{\text{small}} = 1.01 \pm 0.06$ (95% confidence interval), in agreement with the freeze-out picture, suggesting a meaningful β^{eff} , which implies final classical Boltzmann distributions, in this regime. Only a small fraction of the instances, however, correspond to a negligible $Q(s^*)$.

For the majority of instances, $Q(s^*) > 1$, and thus contradicts our working hypothesis that most instances fall in the regime in which one would expect a well-defined freeze-out point and a final classical Boltzmann distribution. working assumption that the instances thermalize according to H_p . The ratio R_Q over the the entire data set is $R_Q \approx 1.20$, substantially higher than the ‘ideal’ $R_Q = 1$. Compared to the rest of the instances, the small $Q(s^*)$ problems are typically easier to solve and are disproportionately smaller in problem size (see e.g. Fig. 16 in Appendix D). The freeze-out picture can also explain the lower-than-ideal effective inverse-temperature ratio $R_\beta = 1.11$ (and higher $R_Q \approx 1.20$). The existence of significant quantum fluctuations outside $A(s) \ll B(s)$ leads to an overestimation of thermal fluctuations in both machines, i.e., to higher effective temperatures, as we indeed observe.

C. High variability in inverse temperature estimates

The magnitude of the error bars on the effective inverse temperatures per instance shown in Fig. 4 (top) reflect the large fluctuations in success probabilities between programming cycles. We discuss various factors that contribute to that variance.

It is known that the location of the freeze-out point (and hence the success probability) has a weak logarithmic dependence on the annealing time [1, 6], with longer anneal times having later freeze-out points because there is more time for fluctuations to take place. We indeed find such an effect (see Figs. 11, 12 of Appendix D), though our results show that this typically accounts for less than a 1% variability between different anneal times and therefore does not explain the spread we observe. If the variation were due purely to statistical variations from cycle to cycle, one would expect statistical fluctuations in success probability P_0 on the

order of $\delta P_0 = \sqrt{P_0(1 - P_0)/N_{\text{anneals}}}$. Fig. 5 (top) shows $R_{\Delta/\delta} = \Delta P_0/\delta P_0$, the ratio of typical magnitude of actual fluctuations in success probabilities ΔP_0 to the expected magnitude of purely statistical fluctuations δP_0 . We find that only around 20% of the instances exhibit fluctuations of success probability $R_{\Delta/\delta}$ below 1. For most instances, typical fluctuations are about an order of magnitude greater than statistical fluctuations, with some fluctuations being considerably greater. We attribute these large ratios, to J -chaos [6] from ICE and other noise, which affect the local fields and coupling parameters within and between cycles. Noise unrelated to programming parameters may also play a role.

Figure 5 (bottom) shows, as a function of problem size, the average variation in β^{eff} , as measured by the ratio of the 95th to 5th percentile values found over all programming cycles. The larger the problem size, the greater the size of the fluctuations. This trend is expected as larger problems, with more couplings, have more potential to be adversely affected by control errors, and other sources of noise [24]. It is critical to understand why these fluctuations scale with problem size, and their root cause, so as to devise strategies to keep these errors from becoming unmanageable as chip sizes increase. For a fixed problem size, we do not observe a clear correlation between success probability and the variance in the β^{eff} estimates (Fig. 6), providing evidence that the fluctuations we observe in Fig. 5 (bottom) are indeed due to differences in problem size and not problem difficulty (though of course the two are related) [25].

IV. CONCLUSIONS

By conducting parallel experiments on two quantum annealers, each operating at a different temperature, we studied key mechanisms determining their output distributions. In particular, we tested the freeze-out conjecture [1, 2, 9] by comparing the performance of the two machines on certain Ising problems, making use of a recent method to accurately estimate the degeneracies of such problems. With a working hypothesis that the output distribution is indeed a Boltzmann distribution of the classical problem Hamiltonian, we calculated the effective inverse temperatures for each instance and machine, β^{eff} , from which we calculated the freeze-out point.

For instances which our results show exhibit negligible quantum fluctuations (small Q), we find a well defined temperature-independent (i.e., machine-independent) freeze-out point, in agreement with the prediction of the freeze-out hypothesis for the small Q regime. This agreement suggests for these instances the output distribution is indeed a classical Boltzmann distribution for H_p , with well defined effective temperature.

Our results also show, however, that for the majority of instances, the estimated freeze-out point is not in the regime of negligible quantum fluctuations, and therefore does not have a well-defined effective temperature, nor

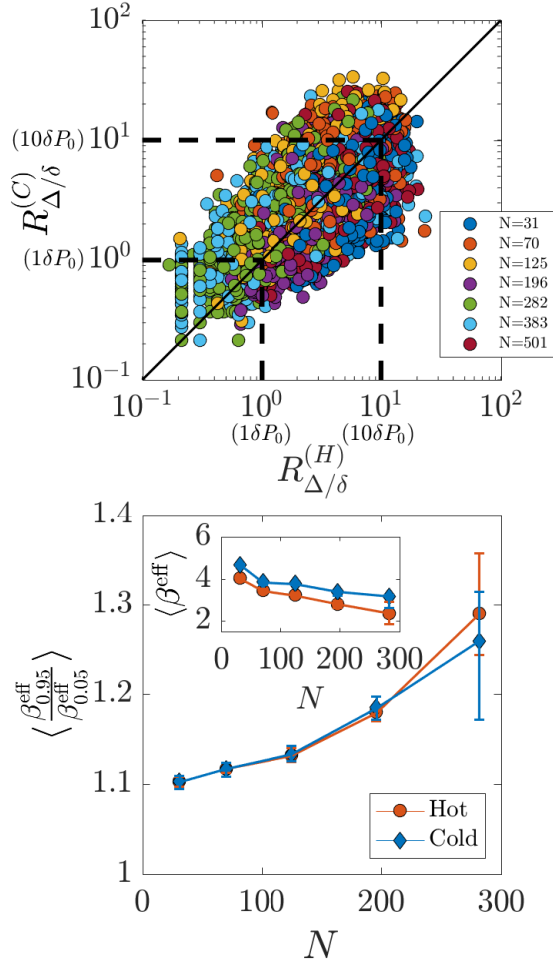


FIG. 5. **Top: Ratio of magnitude of actual fluctuations in success probabilities to magnitude of statistical fluctuations, $R_{\Delta/\delta}$, for the various instances on the hot and cold processors.** For most instances, the fluctuations in success probabilities over programming cycles is far greater (by an order of magnitude) than what one might expect from fluctuations of a purely statistical nature. **Bottom: Typical spread of effective inverse temperature as a function of problem size.** Our measure for spread is the 95th to 5th percentile mean ratio of β^{eff} averaged over instances of each problem size. We take the ratio to overcome any bias from the cold chip recording higher values of β^{eff} (see inset). We find both devices follow a nearly identical trend: fluctuations increase and β^{eff} decreases with problem size. Inset: Median β^{eff} for each problem size. 95% confidence interval error bars obtained by bootstrapping over the instances of each particular N .

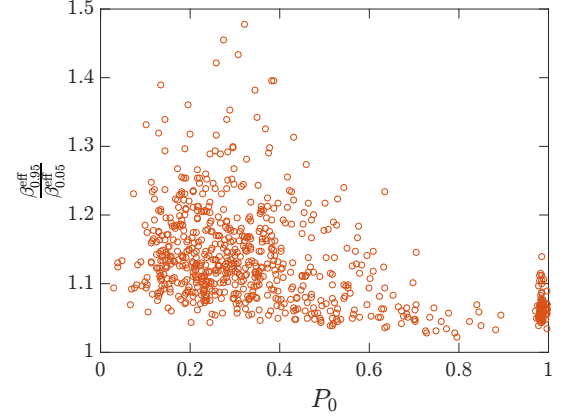


FIG. 6. **Variance with success probability.** Spread of effective inverse temperature as a function of success probability, as determined by the 95th to 5th percentile ratio of β^{eff} over all programming cycles, for $N=70$ (we pick this problem size as it is the one for which we have the most number of instances with reliable data). Data from the hot machine.

is there any reason to believe the output should follow a classical Boltzmann distribution.

Moreover, we also observed that the effective temperatures at different programming cycles can wildly fluctuate. Our data indicates that this effect worsens with larger problem size. These observations show that for future quantum annealers to be effective as Boltzmann samplers, designers must take into account these results, and find ways to ensure that instances thermalize in the $A(s) \ll B(s)$ regime, and such that the effective temperatures are more stable. Moving forward, it would therefore be worthwhile to have additional estimators of temperature, and more robust ways to reconstruct the Boltzmann distribution (i.e., the one we conjecture for small Q) [26].

Promising directions include reducing sources of noise that contribute to intrinsic control errors (ICE) in quantum annealing hardware, and exploring alternate annealing schedules and non-standard drivers to enable more instances to equilibrate at a unique point late enough in the anneal that the quantum fluctuations are negligible. For machine learning, another approach is possible. It is not clear how accurately one needs to sample from Boltzmann distributions for machine learning, or even that Boltzmann distributions are optimal for this purpose. A tantalizing research direction is the use of distributions that have a large quantum component [2], particularly given that certain distributions generated by quantum Hamiltonians are believed to have no efficient classical sampling mechanism [27, 28]. A deeper understanding of these processes will have profound implications for the design of future annealers and the prospects of utilizing quantum annealers as efficient Boltzmann samplers for machine learning and beyond.

ACKNOWLEDGMENTS

We thank Tameem Albash, Mohammad Amin, Salvatore Mandrà and Walter Vinci for useful discussions. We thank the anonymous referees for comments that significantly improved the clarity of the manuscript. Part of the computing resources were provided by the USC Center for High Performance Computing and Communications and the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725. ER would like to acknowledge support from the NASA Advanced Exploration Systems program and NASA Ames Research

Center. Her contributions to this work were also supported in part by the AFRL Information Directorate under grant F4HBKC4162G001 and the Office of the Director of National Intelligence (ODNI) and the Intelligence Advanced Research Projects Activity (IARPA), via IAA 145483. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of ODNI, IARPA, AFRL, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purpose notwithstanding any copyright annotation thereon.

-
- [1] M. H. Amin, Searching for quantum speedup in quasisstatic quantum annealers, *Phys. Rev. A* **92**, 052323 (2015).
 - [2] M. H. Amin, E. Andriyash, J. Rolfe, B. Kulchytskyy, R. Melko, Quantum Boltzmann Machine, ArXiv e-prints (2016), [arXiv:1601.02036 \[quant-ph\]](#).
 - [3] M. Benedetti, J. Realpe-Gómez, R. Biswas, and A. Perdomo-Ortiz, Estimation of effective temperatures in quantum annealers for sampling applications: A case study with possible applications in deep learning, *Phys. Rev. A* **94**, 022308 (2016).
 - [4] B. H. Zhang, G. Wagenbreth, V. Martin-Mayor, and I. Hen, Advantages of Unfair Quantum Ground-State Sampling, *Scientific Reports* **7**, 1044 (2017).
 - [5] H. G. Katzgraber, F. Hamze, Z. Zhu, A. J. Ochoa, and H. Muñoz-Bauza, Seeking Quantum Speedup Through Spin Glasses: The Good, the Bad, and the Ugly, *Phys. Rev. X* **5**, 031026 (2015).
 - [6] V. Martin-Mayor and I. Hen, Unraveling Quantum Annealers using Classical Hardness, *Scientific Reports* **5**, 15324 (2015).
 - [7] J. Marshall, V. Martin-Mayor, and I. Hen, Practical engineering of hard spin-glass instances, *Phys. Rev. A* **94**, 012320 (2016).
 - [8] S. H. Adachi and M. P. Henderson, Application of Quantum Annealing to Training of Deep Neural Networks, ArXiv e-prints (2015), [arXiv:1510.06356 \[quant-ph\]](#).
 - [9] M. W. Johnson *et al.*, Quantum annealing with manufactured spins, *Nature* **473**, 194 (2011).
 - [10] A. D. King and C. C. McGeoch, Algorithm engineering for a quantum annealing platform, ArXiv e-prints (2014), [arXiv:1410.2628 \[quant-ph\]](#).
 - [11] Z. Zhu, A. J. Ochoa, S. Schnabel, F. Hamze, and H. G. Katzgraber, Best-case performance of quantum annealers on native spin-glass benchmarks: How chaos can affect success probabilities, *Phys. Rev. A* **93**, 012317 (2016).
 - [12] S. Boixo, V. N. Smelyanskiy, A. Shabani, S. V. Isakov, M. Dykman, V. S. Denchev, M. H. Amin, A. Y. Smirnov, M. Mohseni, and H. Neven, Computational multiqubit tunnelling in programmable quantum annealers, *Nature Communications* **7**, 10327 EP (2016).
 - [13] D. Venturelli, S. Mandrà, S. Knysh, B. O’Gorman, R. Biswas, and V. Smelyanskiy, Quantum Optimization of Fully Connected Spin Glasses, *Phys. Rev. X* **5**, 031040 (2015).
 - [14] M. Nifle and H. J. Hilhorst, New critical-point exponent and new scaling laws for short-range Ising spin glasses, *Phys. Rev. Lett.* **68**, 2992 (1992).
 - [15] M. Ney-Nifle, Chaos and universality in a four-dimensional spin glass, *Phys. Rev. B* **57**, 492 (1998).
 - [16] F. Krzakala and J. P. Bouchaud, Disorder chaos in spin glasses, *Europhys. Lett.* **72**, 472 (2005).
 - [17] H. G. Katzgraber and F. Krzakala, Temperature and Disorder Chaos in Three-Dimensional Ising Spin Glasses, *Phys. Rev. Lett.* **98**, 017201 (2007).
 - [18] T. Albash, V. Martin-Mayor, and I. Hen, Temperature Scaling Law for Quantum Annealing Optimizers, *Phys. Rev. Lett.* **119**, 110502 (2017).
 - [19] The term ‘effective temperature’ is somewhat of a misuse as it may imply thermalization of the system whereas in fact it may not be the case.
 - [20] One machine is owned by Lockheed-Martin, housed at USC’s Information Sciences Institute and the other, by a NASA-USRA-Google collaboration and housed inside the NASA Ames Research Center.
 - [21] I. Hen, J. Job, T. Albash, T. F. Rønnow, M. Troyer, and D. A. Lidar, Probing for quantum speedup in spin-glass problems with planted solutions, *Phys. Rev. A* **92**, 042325 (2015).
 - [22] F. Wang and D. P. Landau, Determining the density of states for classical statistical models: A random walk algorithm to produce a flat histogram, *Phys. Rev. E* **64**, 056101 (2001).
 - [23] If we calculate the ratio via the (least squares) gradient of Fig. 4 (top), we find it is $R = 1.14$, also far below the physical ratio.
 - [24] The increase in fluctuations with problem size we observe in Fig. 5 (bottom) is most likely an underestimate of the full effect. Since our criterion for discarding instances is convergence of the WL algorithm, those instances that do not appear in the figure exhibit fluctuations of larger magnitudes, as there is a known strong positive correlation between WL convergence, i.e., its classical hardness, and J -chaos (see, e.g., Ref. [6]).
 - [25] We also discount other minor effects, such as the known logarithmic dependence on anneal time, in Appendix D.
 - [26] J. Raymond, S. Yarkoni, and E. Andriyash, Global Warming: Temperature Estimation in Annealers, *Front-*

- tiers in ICT **3**, 23 (2016).
- [27] C. Semay and L. Ducobu, Quantum and classical probability distributions for arbitrary Hamiltonians, *European Journal of Physics* **37**, 045403 (2016).
 - [28] R. W. Robinett, Quantum and classical probability distributions for position and momentum, *Am. J. Phys.* **63**, 823 (1995).
 - [29] P. I. Bunyk, E. M. Hoskinson, M. W. Johnson, E. Tolkacheva, F. Altomare, A. Berkley, R. Harris, J. P. Hilton, T. Lanting, A. Przybysz, and J. Whittaker, Architectural Considerations in the Design of a Superconducting Quantum Annealing Processor, *Applied Superconductivity, IEEE Transactions on* **24**, 1 (Aug. 2014).
 - [30] V. Choi, Minor-embedding in adiabatic quantum computation: II. Minor-universal graph design, *Quant. Inf. Proc.* **10**, 343 (2011).
 - [31] W. Barthel, A. K. Hartmann, M. Leone, F. Ricci-Tersenghi, M. Weigt, and R. Zecchina, Hiding Solutions in Random Satisfiability Problems: A Statistical Mechanics Approach, *Phys. Rev. Lett.* **88**, 188701 (2002).
 - [32] F. Krzakala and L. Zdeborová, Hiding Quiet Solutions in Random Constraint Satisfaction Problems, *Phys. Rev. Lett.* **102**, 238701 (2009).
 - [33] S. Boixo, T. F. Rønnow, S. V. Isakov, Z. Wang, D. Wecker, D. A. Lidar, J. M. Martinis, and M. Troyer, Evidence for quantum annealing with more than one hundred qubits, *Nat. Phys.* **10**, 218 (2014).
 - [34] T. F. Rønnow, Z. Wang, J. Job, S. Boixo, S. V. Isakov, D. Wecker, J. M. Martinis, D. A. Lidar, and M. Troyer, Defining and detecting quantum speedup, *Science* **345**, 420 (2014).
 - [35] M. Mezard, G. Parisi and M.A. Virasoro, *Spin Glass Theory and Beyond*, World Scientific Lecture Notes in Physics (World Scientific, Singapore, 1987).

Appendix A: The Google-NASA-USRA (‘cold’) and Lockheed-Martin-USC (‘hot’) D-Wave Two processors

The quantum annealer used in our work is the D-Wave Two (DW2) device [29]. This device is designed to solve optimization problems by evolving a known initial configuration — the ground state of a transverse field $H_d = -\sum_i \sigma_i^x$, where σ_i^x is the Pauli spin-1/2 matrix acting on spin i — towards the ground state of a classical Ising-model Hamiltonian which serves as a cost function that encodes the problem that is to be solved:

$$H_p = \sum_{\langle i,j \rangle} J_{ij} \sigma_i^z \sigma_j^z + \sum_i h_i \sigma_i^z. \quad (\text{A1})$$

The variables $\{\sigma_i^z\}$ denote either classical Ising-spin variables that take values ± 1 or Pauli spin-1/2 matrices, the $\{J_{ij}\}$ are programmable coupling parameters, and the $\{h_i\}$ are programmable local longitudinal fields. The N spin variables are realized as superconducting flux qubits and occupy the vertices of the D-Wave ‘Chimera’ hardware graph [29, 30]. Here, $\langle i, j \rangle$ denotes summation over the edges of the graph. The union of the two D-Wave Chimera graphs is given in Fig. 2 – this is the graph all of our problem instances were defined on.

These machines evolve the full Hamiltonian via

$$H(s) = A(s)H_d + B(s)H_p. \quad (\text{A2})$$

The way in which the strength of the initial (H_d) and final (H_p) Hamiltonians evolve is given by the parameters $A(s)$ and $B(s)$, where $s = t/\mathcal{T} \in [0, 1]$ is the annealing time. Here, \mathcal{T} is the total annealing time, ranging between 20 μ s and 20ms on these devices. The annealing schedule is given in Fig. 7 for each machine.

In Fig. 8 we show the temperature log of the D-Wave chips during the time which we collected our data.

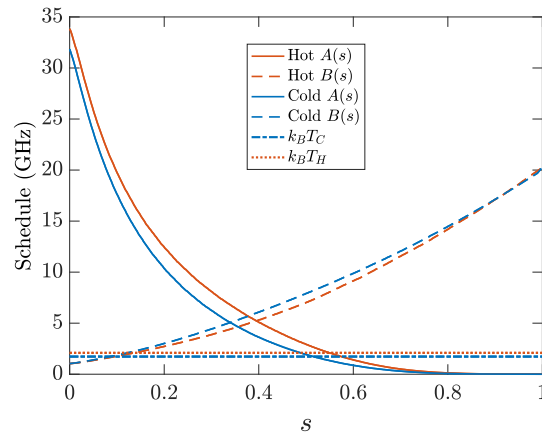


FIG. 7. **Annealing schedules of the USC (hot) and NASA (cold) DW2 processors.** Annealing schedule [see Eq. (A2)] in GHz as a function of dimensionless annealing time $s = t/\mathcal{T}$. We also plot the temperatures ($\hbar = 1$) of the devices (see legend).

Appendix B: Generation of instances

For the generation of instances, we have chosen in this work to study problems constructed around ‘planted solutions’ – an idea borrowed from constraint satisfaction (SAT) problems [31, 32]. In these problems, the planted solution represents a ground state configuration of Eq. (A1) that minimizes the energy and is known in advance. This knowledge circumvents the need to verify the ground state energy using exact (provable) solvers, which rapidly become too expensive computationally as the number of variables grows, and which were employed in earlier benchmarking studies [33, 34]. Moreover, these problems are known to possess different degrees of ‘tunable hardness’, achieved by adjusting the amount of frustration (see Ref. [35]) which we will use. Last, studying this type of problems allows us to devise an algorithm to find all minimizing configurations of the generated instances. Knowledge of both energy levels as their degeneracies is essential for the calculation of effective temperatures of the instances.

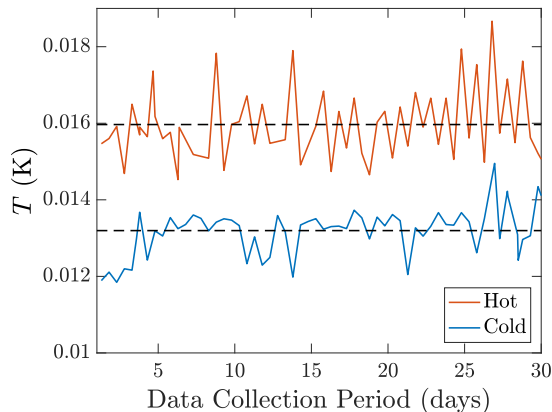


FIG. 8. **Temperature logs.** Temperature log for the USC (hot) and NASA (cold) machines during the period of time we performed our experiments. The dashed black lines represent the mean of the data sets. The mean temperature of the ‘hot’ USC machine was about $T_H = 16.0$ mK, and the mean temperature of the ‘cold’ NASA machine was about $T_C = 13.2$ mK, with ratio $T_H/T_C \approx 1.21$. Note, the temperature data is sparse, sampled only twice per 24 hours.

To generate the instances for the experiment we follow the guidelines introduced in Ref. [21]. We generate 13 groups of 100 instances, for each of 7 different sub-Chimera sizes with $L = 2 \dots 8$, see Fig. 2 (i.e., 9100 total instances). These 13 groups differ in the ratio α of number of clauses (or loops) to number of qubits contained in each instance. For every fixed ratio α the range of integer-valued J_{ij} ’s, which we denote by J_{\max} is fixed across the different problem sizes. D-Wave further rescales all coupling values such that the encoded values, $\tilde{J} \in [-1, 1]$ which implies that for every fixed α both the range of J -values as well as the spacings between them is identical across different problem sizes.

Appendix C: Wang-Landau entropic sampler

As explained in the main text, we employed a Wang-Landau entropic sampler to estimate the degeneracy of the energy levels for our generated planted-solution instances. This algorithm performs essentially a random walk over the energy landscape, where updates at each step in the algorithm are such that an approximately flat histogram of visited energies is produced. We follow the same methodology as originally described in [22]. Our histogram was considered ‘flat’ when the lowest sampled energy level has been visited at least 80% of the mean of the entire histogram.

We performed 20 independent Wang-Landau runs, each up to 10^9 steps for each of our instances. We then averaged over these 20 runs which provided our estimate of degeneracies for each instance. We then discarded

any instances for which the ground or first excited state degeneracies did not match that for the exact solution counter (up to 5% error). This meant we had accurate degeneracy data for problems up to 282 qubits in size.

Appendix D: Effects or lack thereof of additional experimental parameters

1. Success probabilities

In Fig. 9 we show the histogram of the success probabilities for the two machines, for all of the instances. We see the cold (NASA) machine clearly outperforms the hot (USC) machine—we expect, due to the colder operating temperature.

2. Programming cycles

In Fig. 10 we compare two different programming cycles (from different days) on the same machine, showing consistency over different runs.

3. Anneal times

We also study the effect of varying anneal time on success probability in Figs. 11 and 12. We see that there is only a very weak (logarithmic) dependence on anneal time, in accordance with [1, 6], and moreover, it is seemingly not correlated with problem size.

4. Ratio of number of clauses to number of qubits

Fig. 13 shows that (within error bars) our computation of the inverse temperature ratio between the two machines is unaffected by changing the ratio of number of clauses (or loops) to number of qubits α , and moreover, that the thermal ratio is well above the measured ratio (far outside of the error bars mostly).

Appendix E: Additional data

Relating to Fig. 4 of the main text, we produce a ‘heat map’, Fig. 14, for β^{eff} , and for Q^* (defined in main text), for each machine, which shows the number of instances found in each small region. We notice that in the upper figure that the larger the effective inverse temperature, more instances fall within the ‘thermal region’, indicating a stronger dependence on the temperature for these instances. These instances correspond to the ones which freezeout at a later point in the anneal, and thus a smaller Q^* in the lower figure. In this lower figure, we observe the fit is closer to the ‘ideal’ $y = x$ for smaller Q^* , and it

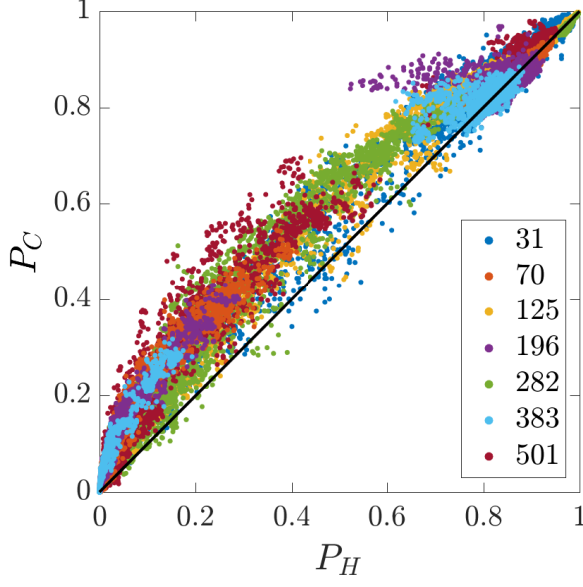


FIG. 9. **Success probability.** Probability of success (how often the ground state energy is correctly identified) of all instances and programming cycles on the two machines ('hot' USC machine P_H , and the 'cold' NASA machine P_C). Each point is a randomly chosen programming cycle (for the same instance on each machine). Number of qubits given by legend.

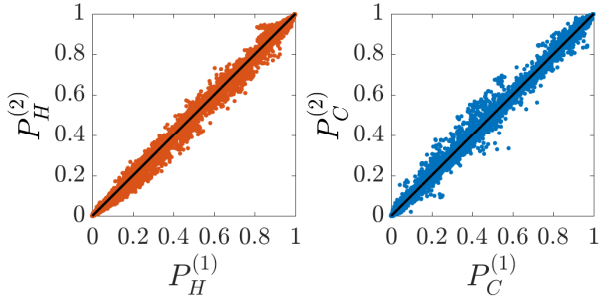


FIG. 10. **Machine correlation.** We compare the results of two programming cycles for each instance on each machine. We see the data aligns nicely along $y = x$, albeit with sizable fluctuations (as one would expect). Compare this with Fig. 9, where data clearly deviates from $y = x$. The 'hot' USC machine is on the left, and the 'cold' NASA machine on the right.

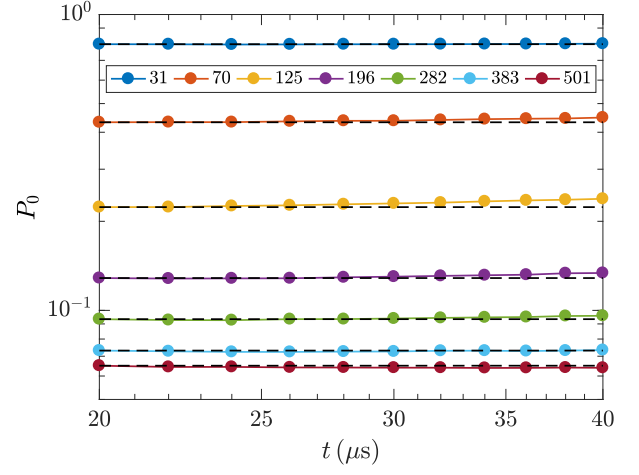


FIG. 11. **D-Wave success with anneal time.** Average probability of success, P_0 (for the hotter USC machine), against anneal time (log scale), for different problem sizes (see legend). Each point averaged over two programming cycles, of $N_{\text{anneals}} = 20,000$ anneals each. The black dash lines correspond to the value of P_0 for $t = 20 \mu\text{s}$. One can see in general a slight increase in success probability with t .

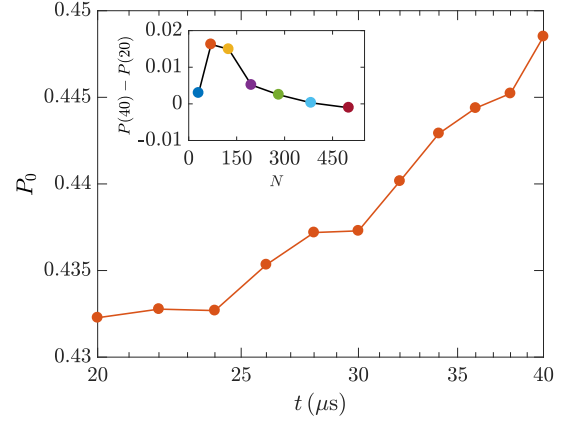


FIG. 12. **D-Wave success with anneal time, for single problem size.** Average probability of success, P_0 (for the hotter USC machine), against anneal time (log scale), for problem size 70 ($L = 3$). Note the approximate linear relationship. Inset: Difference in success probability between $t = 40$ and $t = 20 \mu\text{s}$ as measured by $P(40) - P(20)$, where $P(t)$ is defined as P_0 for anneal time t .

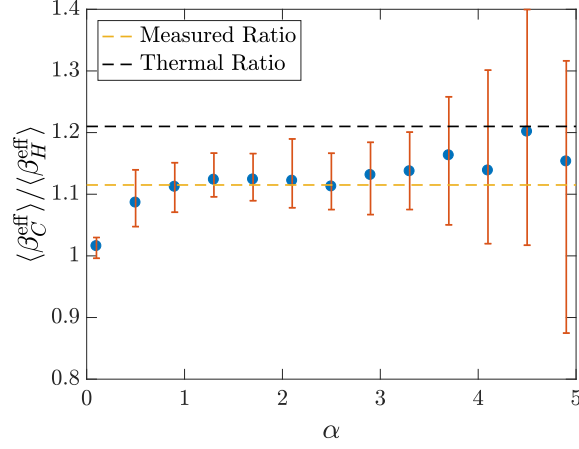


FIG. 13. **Effective inverse temperature ratio against clause density.** As explained in the text, we generated many problems for each graph size $L = 2, \dots, 8$. We picked fixed and identical values of α (the ratio of the number of loops to the number of qubits) for each qubit size. We find that α has (within the 95% confidence interval error bars obtained from bootstrapping the data) mostly no effect on the effective temperature ratio we measured from our data (1.11, calculated as the median over all instances), and that it is much below the thermal ratio value of 1.21. Note the smallest $\alpha = 0.1$ instances are very easy to solve ($P_0 \approx 1$), making it hard to distinguish performance (hence temperature) differences between the machines.

deviates above this line for larger Q^* (we discuss in more detail in the main text).

Also relating to Fig. 4 of the main text, we produce Figs. 15, 16, which plot β^{eff} , and Q^* , for each machine, but split up by problem size. One can see that typically the larger problems exhibit lower values of β^{eff} , and likewise, larger values of Q^* , indicating these are in fact not thermalizing according to a Boltzmann distribution.

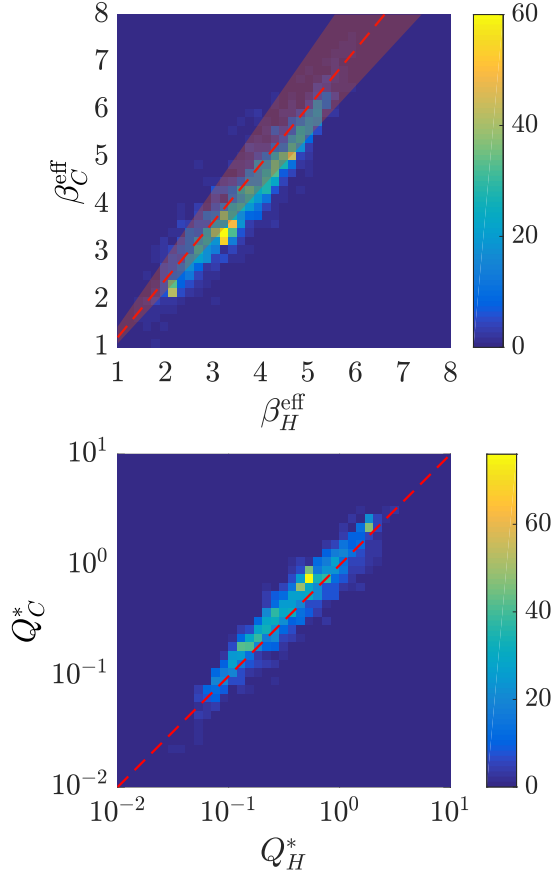


FIG. 14. **Density plots for β^{eff} and Q^* .** The corresponding ‘heat map’ for Fig. 2 of the main text. Color indicates the number of points in each region (given by color bar on the right hand side). Top: Red dash line is the ‘ideal’ thermal ratio (i.e. if the ratio of the effective inverse temperatures were the same as the physical ‘thermal’ inverse temperatures). The variance in physical temperature fluctuations is given by the red semi-transparent region. Bottom: Red dash line is $y = x$.

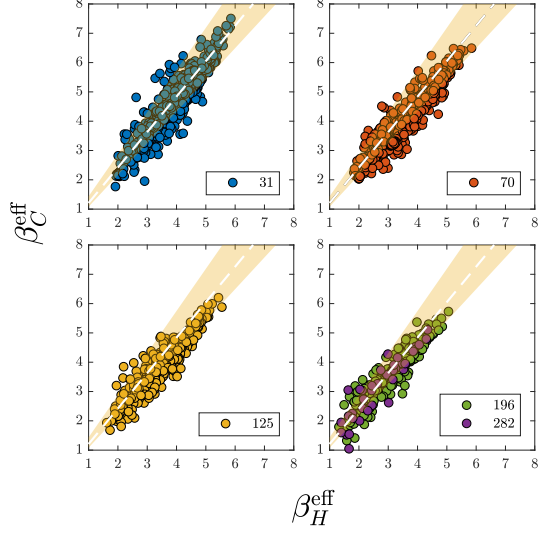


FIG. 15. β^{eff} for different problem sizes. These plots show the effective inverse temperature found by each machine (and for each instance), for the 5 different problem sizes (see legend) for which we have reliable degeneracy data. White dash line is the ‘ideal’ thermal ratio (i.e. if the ratio of the effective inverse temperatures were the same as the physical ‘thermal’ inverse temperatures). The variance in physical temperature fluctuations is given by the yellow semi-transparent region.

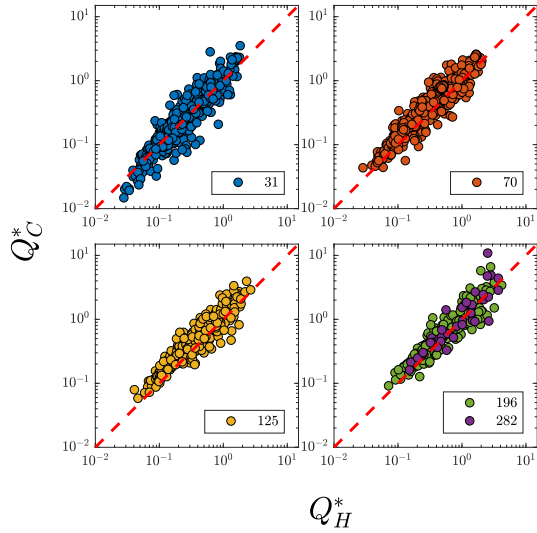


FIG. 16. Q^* for different problem sizes. These plots show the median Q^* found by each machine (and for each instance), for the 5 different problem sizes (see legend) for which we have reliable degeneracy data. Red dash line is $y = x$.