# Stochastic Neuromorphic Circuits for Solving MAXCUT

Bradley H. Theilman
*Neural Exploration and Research Lab*
*Sandia National Laboratories*
Albuquerque, New Mexico
bhtheil@sandia.gov

Yipu Wang
*Discrete Math and Optimization*
*Sandia National Laboratories*
Albuquerque, New Mexico
yipwang@sandia.gov

Ojas Parekh
*Discrete Math and Optimization*
*Sandia National Laboratories*
Albuquerque, New Mexico
odparek@sandia.gov

William Severa
*Neural Exploration and Research Lab*
*Sandia National Laboratories*
Albuquerque, New Mexico
wmsever@sandia.gov

J. Darby Smith
*Neural Exploration and Research Lab*
*Sandia National Laboratories*
Albuquerque, New Mexico
jsmit16@sandia.gov

James B. Aimone
*Neural Exploration and Research Lab*
*Sandia National Laboratories*
Albuquerque, New Mexico
jbaimon@sandia.gov

*Abstract*—Finding the maximum cut of a graph (MAXCUT) is a classic optimization problem that has motivated parallel algorithm development. While approximate algorithms to MAXCUT offer attractive theoretical guarantees and demonstrate compelling empirical performance, such approximation approaches can shift the dominant computational cost to the stochastic sampling operations. Neuromorphic computing, which uses the organizing principles of the nervous system to inspire new parallel computing architectures, offers a possible solution. One ubiquitous feature of natural brains is stochasticity: the individual elements of biological neural networks possess an intrinsic randomness that serves as a resource enabling their unique computational capacities. By designing circuits and algorithms that make use of randomness similarly to natural brains, we hypothesize that the intrinsic randomness in microelectronics devices could be turned into a valuable component of a neuromorphic architecture enabling more efficient computations. Here, we present neuromorphic circuits that transform the stochastic behavior of a pool of random devices into useful correlations that drive stochastic solutions to MAXCUT. We show that these circuits perform favorably in comparison to software solvers and argue that this neuromorphic hardware implementation provides a path for scaling advantages. This work demonstrates the utility of combining neuromorphic principles with intrinsic randomness as a computational resource for new computational architectures.

## I. INTRODUCTION

Despite the heavy requirements for noise-free operation placed on the components of conventional computers, random numbers play a crucially important role in many parallel computing problems arising in different scientific domains. Because current random number generation occurs largely in software, the required randomness in these systems is plagued by the same memory-processing bottlenecks that limit ordinary computation. Current work in material science and microelectronics is demonstrating the feasibility of constructing stochastic microelectronic devices with controllable statistics for probabilistic neural computing [22]. These devices show scalability properties that forecast the ability to generate random numbers in-situ with the processing elements, bypassing this bottleneck.

Stochasticity is an inherent property of physical systems, both natural and artificial. Physical computers are able to approximate ideal computations because much effort has been expended in developing electronic technology that minimizes the influence of universal electronic "noise." As microelectronics get smaller and the scale of our computations get larger, current computational paradigms require even more stringent limits on the influence of this electronic noise, and these limits become severe constraints on the scalability of existing computational architectures.

In contrast, natural brains are examples of highly parallel computational systems that achieve amazingly efficient computational performance in the face of ubiquitous noise. There are on the order of $10^{15}$

synapses in a human brain, and each one is stochastic: its probability of successfully transmitting a signal to a downstream neuron ranges from 0.1 to 0.9 [20]. Each synapse is activated about once per second on average, so, the brain generates about $10^{15}$ random numbers per second [22]. Compare this to the reliability of transistor switching in conventional computers, where the probability of failure is less than $10^{-14}$ [20]. It is unknown precisely how brains deal with this stochasticity, but its pervasiveness strongly suggests that the brain uses its own randomness as a computational resource rather than treating it as a defect that must be eliminated. This suggests that a new class of parallel computing architectures could emerge from combining the computational principles of natural brains with physical sources of intrinsic randomness. This would allow the natural stochasticity of electronic devices to play a part in large-scale parallel computations, relieving the burden imposed by requiring absolute reliability.

Realizing the potential of probabilistic neural computation requires rethinking conventional parallel algorithms to incorporate stochastic elements from the bottom up. Additionally, techniques for controlling the randomness must be developed so that useful random numbers can be produced efficiently from the desired distributions. In this work, we propose neuromorphic circuits that demonstrate the capacity for intrinsic randomness to solve parallel computing problems and techniques for controlling device randomness to produce useful random numbers.

MAXCUT is a well known, NP-complete problem that has practical applications and serves as a model problem and testbed for both classical and beyond-Moore algorithm development [4], [9], [15], [19]. The problem requires partitioning the vertices of a graph into two disjoint classes such that the number of edges that span classes is maximized. MAXCUT has several stochastic approximation algorithms, which makes it an ideal target for developing new architectures leveraging large-scale parallel stochastic circuit elements for computational benefit.

Stochastic approximation algorithms are compared via their approximation ratio, which is the ratio of the expected value of a stochastically generated solution to the maximum possible value. The stochastic approximation to MAXCUT with the largest known approximation ratio is the Goemans-Williamson algorithm [9]. The Goemans-Williamson algorithm provides the best approximation ratio achievable by any polynomial-time algorithm under the Unique Games Conjecture [19]. To generate solutions, this algorithm requires sampling from a Gaussian distribution with a specific covariance matrix obtained by solving a semi-definite program related to the adjacency matrix of the graph. Our first neural circuit implements this sampling step

by using simple neuron models to transform uniform device randomness into the required distribution. This demonstrates the use of neuromorphic principles to transform an intrinsic source of randomness into a computationally useful distribution.

Another stochastic approximation for MAXCUT is the Trevisan algorithm [27], [29]. Despite having a worse theoretical approximation ratio, in practice this algorithm generates solutions on par with the Goemans-Williamson algorithm [21]. To generate solutions, this algorithm requires computing the minimum eigenvector of the normalized adjacency matrix. Our second neuromorphic circuit implements this algorithm using the same circuit motif as above to generate random numbers with a specific correlation, but instead of sampling cuts from this distribution, we use these numbers to drive a synaptic plasticity rule (Oja's rule) inspired by the Hebbian principle in neuroscience [24]. This learning rule can be shown to converge to the desired eigenvector, from which the solution can be sampled. This circuit solves the MAXCUT problem entirely within the circuit, without requiring any external preprocessing, demonstrating the capacity of neuromorphic circuits driven by intrinsic randomness to solve parallel computationally-relevant problems.

Neuromorphic computing is having increasing impacts on non-cognitive problems relevant for parallel computing [1]. Unlike other hardware approaches to MAXCUT, our contributions directly instantiate state-of-the-art MAXCUT approximation algorithms on arbitrary graphs without requiring costly reconfiguration or conversion of the problem to an Ising model with pairwise interactions [10], [11], [30]. Our use of hardware resources is scalable, requiring one neuron and one random device per vertex, and thus more efficient than parallel implementations of MAXCUT using GPUs [8]. These properties make our contributions valuable to the expanding field of beyond-Moore parallel algorithms.

## II. MAXCUT ALGORITHMS

### A. The Goemans-Williamson MAXCUT Algorithm

Given an $n$-vertex, $m$-edge graph $\mathcal{G} = (V, E)$ with vertex set $V$ and edge set $E$, the MAXCUT problem seeks a partition of the vertices into two disjoint subsets, $V = V_{-1} \cup V_1$ such that the number of edges that cross between the two subsets is maximized. By assigning either of the values $\{-1, 1\}$ to each vertex, the MAXCUT problem is equivalent to maximizing the function

$$\max_{v} \quad \frac{1}{2} \sum_{ij \in V} A_{ij}(1 - v_i v_j)$$
$$\text{s.t.} \quad v \in \{-1, 1\}^n.$$

Here, $A_{ij}$ is the adjacency matrix of the graph $\mathcal{G}$. Let OPT($\mathcal{G}$) be the maximum value of this function.

MAXCUT is known to be NP-complete [15]. Goemans and Williamson [9] described a relaxation of the above integer programming problem that yields an efficient approximation to MAXCUT with an approximation ratio of 0.878. The relaxation replaces the integer programming problem with a semidefinite programming problem given by

$$\max_{w} \quad \frac{1}{2} \sum_{ij} A_{ij}(1 - w_i \cdot w_j)$$
$$\text{s.t.} \quad w_i \in \mathcal{S}^{n-1},$$

where $\mathcal{S}^{n-1}$ is the $(n-1)$-dimensional unit sphere in $\mathbb{R}^n$. Let SDP($\mathcal{G}$) be value of the optimal solution of this semidefinite programming problem. Note that OPT($\mathcal{G}$) $\leq$ SDP($\mathcal{G}$).

The solution of this semidefinite programming problem is a set of unit vectors $w_i$, one for each vertex in the graph. Given these vectors, a graph cut is generated by taking a random hyperplane through the origin and assigning the value $+1$ to vertices with vectors above the plane and $-1$ to vertices with vectors below the plane.

One can see that the Goemans-Williamson algorithm has two steps: in the first step we solve an SDP, and in the second we round each unit vector $w_i$ to an integer $z_i \in \{-1, +1\}$, where $i \in V$. Bertsimas and Ye [6] observed that the rounding step can be implemented by sampling dependent standard normal random variables, with one variable per vertex. Specifically, suppose for each vertex $i$ we have a random variable $X_i$ following the standard normal distribution, and furthermore for each pair of vertices $i$ and $j$ the covariance between $X_i$ and $X_j$ is $w_i \cdot w_j$, where $w_i$ and $w_j$ are the unit vectors in the solution to the SDP. One can show that such a set of dependent random variables exists. Now define a (random) cut by assigning $+1$ to vertices $i$ where $X_i$ is positive and assigning $-1$ to vertices $i$ where $X_i$ is negative. One can show that the resulting cut has the same approximation guarantees as the cut returned by the Goemans-Williamson algorithm. Hence in this paper we will sometimes refer to the rounding step as the sampling step.

### B. The Trevisan (Simple Spectral) Algorithm

The Trevisan algorithm is another random approximation algorithm for MAXCUT [29]. Though it has a worse theoretical approximation ratio (0.631) [27] than the Goemans-Williamson algorithm, in practice it can perform just as well and has speed advantages [21]. Here we consider a slight modification of the full Trevisan algorithm we refer to as the Trevisan Simple Spectral algorithm [21].

Given a graph $G = (V, E)$ with adjacency matrix $A$ and diagonal degree matrix $D$, we compute the normalized adjacency matrix $\mathcal{A} = D^{-1/2}AD^{1/2}$. Next, the eigenvector corresponding to the minimum eigenvalue of the matrix $I + \mathcal{A}$ is computed. The graph cut is obtained by thresholding the values of this eigenvector by sign. If $\mathbf{u}$ is the minimum eigenvector of $\mathcal{A}$, then the graph cut is given by

$$v_i = \begin{cases} -1 & \mathbf{u}_i \leq 0 \\ 1 & \mathbf{u}_i > 0 \end{cases}$$

### III. Neuromorphic Concepts

#### A. Stochastic Devices

Physical microelectronics display intrinsic stochasticity due the physics behind their operation. Typically this stochasticity is observed as random switching between two or more states. While normally a nuisance, the details of this stochastic behavior are under active research to develop devices with tunable statistics for probabilistic computing applications [5], [22], [25]. In our work, we idealize stochastic devices as analogous to "coin flips" such that at any given time, the device can be in one of two states ("heads" or "tails"; "0" or "1") with a specific probability. In our circuits, we assume random devices behave as fair coins. That is, each state has a probability of 0.5. Thus, a random device is modeled as a source for a random bit stream with equal probabilities of 0 or 1. Magnetic tunnel junctions [18], [25] and tunnel diodes [5] are examples of two classes of devices actively being developed to meet these requirements.

#### B. Leaky Integrate and Fire Neurons

The leaky integrate and fire (LIF) neuron is a simplified model of biological neurons, readily implemented in hardware [14], that captures a biological neuron's capacity for temporal integration of synaptic inputs along with discontinuous spiking. The model integrates synaptic currents with a membrane capacitance into a membrane potential that is continuously discharged by a leak conductance. When the integrated membrane potential reaches some threshold, a spike is emitted and the membrane potential is reset to some defined value. In between spike events, the membrane potential evolves according to the differential equation

$$C\frac{dV}{dt} = -\frac{V}{R} + I_{\text{tot}}.$$

Here, $V$ is the membrane potential, $C$ and $R$ are the membrane capacitance and leak resistance, respectively, and $I_{\text{tot}}$ is the total synaptic input current.

When a single LIF neuron receives large numbers of stochastic input currents, the membrane potential approximates a one-dimensional random walk [20].

The leak conductance stabilizes this walk around an analytically computable mean

$$\langle V \rangle = R \langle I_{\text{tot}} \rangle$$

and variance

$$\text{Var}(V) = \frac{R}{C}\text{Var}(I_{\text{tot}}).$$

## C. LIF Covariances

For a population of $n$ LIF neurons integrating random binary inputs generated by $r$ random devices, the expression for the membrane potential dynamics of a single LIF neuron becomes

$$C\frac{dV_i}{dt} = -\frac{V_i}{R} + \sum_\alpha W_{i\alpha}s_\alpha.$$

Here $W_{i\alpha}$ is the real-valued connection weight between device $\alpha$ and LIF neuron $i$. The variable $s_\alpha$ is the state of device $\alpha$ and takes the values $\{0,1\}$.

Shared or inverted input between two LIF neurons induces correlations or anticorrelations in their membrane potentials, respectively. The expression for the covariance between the membrane potentials of neurons $i$ and $j$ is

$$\text{Cov}(V_i, V_j) = \frac{R}{C}\sum_{\alpha\beta} W_{i\alpha}W_{j\beta}\text{Cov}(s_\alpha, s_\beta).$$

In other words, the LIF membrane covariances are a linear transformation of the covariances of the random device pool. The device covariance matrix defines an inner product on the space of weight vectors for each LIF neuron. If the devices are independent, then the device covariance matrix is diagonal. Thus, the LIF neuron population transforms the device randomness into a set of Gaussian processes with covariance proportional to the Gram matrix of the weight vectors. In what follows, choosing the weights appropriately allows this circuit motif to supply random samples with the appropriate covariances for the stochastic MAXCUT approximation algorithms.

## D. Synaptic Plasticity: Oja's Rule

In neuroscience, the guiding principle of synaptic plasticity is captured by the adage "neurons that fire together, wire together." This is the Hebbian learning principle [13]. If $\mathbf{w}$ is the weight vector between presynaptic neuron activity $\mathbf{x}$ and postsynaptic neuron activity $y$, the simplest instantiation of this principle is given the formula

$$\Delta\mathbf{w} = y\mathbf{x}$$

As stated, this rule is unstable. Oja presented a modification to this plasticity rule that preserved the Hebbian principle but enforced weight stability [23]. Oja's rule is given by the formula

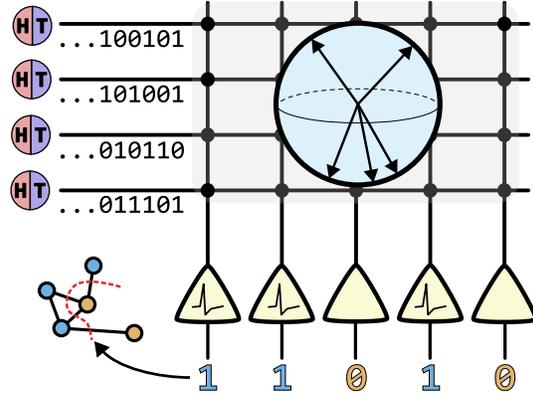$$\Delta\mathbf{w} = y(\mathbf{x} - y\mathbf{w})$$



Figure 1. LIF-Goemans-Williamson circuit implementing the sampling stage of the Goemans-Williamson algorithm. Spikes from the LIF population correspond to binary labels on the vertices of the graph, defining the cut. The covariances of the LIF membrane potentials are determined by the weight matrix from the random devices (left) to the LIF population. Each LIF neuron's weight vector is set proportional to a vector determined through the solution to the Goemans-Williamson SDP.

Oja proved that under mild assumptions this rule forces the weight vector to converge to the first principle component of the covariance matrix of the inputs, or equivalently the eigenvector corresponding the the largest eigenvalue.

By considering anti-Hebbian plasticity, Oja derived a related, stabilized learning rule that converges to the *minimum* eigenvector of the covariance matrix [24]

$$\Delta\mathbf{w} = -y\mathbf{x} + (y^2 + 1 - \mathbf{w}^T\mathbf{w})\mathbf{w}$$

By providing inputs with covariance proportional to the adjacency matrix of the graph as used in Trevisan's algorithm, Oja's anti-Hebbian rule can find the minimum eigenvector of this matrix, yielding an approximate solution to MAXCUT.

## IV. CIRCUITS

### A. LIF-Goemans-Williamson

Figure 1 shows a neural circuit that implements the sampling step of the Goemans-Williamson algorithm. The requirement is to generate binarized samples from a Gaussian distribution with specified covariance matrix $C$. We refer to this circuit with the abbreviation LIF-GW.

For a graph $G$, the Goemans-Williamson SDP is solved to yield a set of $n$ unit vectors in $r$ dimensions, where $r$ is the rank of the solution and $n$ is the number of vertices. These vectors can be combined into the $n$ by $r$ dimensional matrix $W_{GW}$.

The circuit consists of a pool of $r$ random devices connected to $n$ LIF neurons. The synaptic weights between the devices and the neurons are chosen proportional to the corresponding entries in $W_{GW}$. The precise magnitudes of these weights are not critical; what matter are their relative values, as

these ratios determine the LIF covariances. This allows the circuit to be adapted to specific hardware implementations imposing constraints on the range of available weights. For our tests, we used a fixed rank of 4 for all graphs.

Choosing the weights proportional to the solution to the SDP yields membrane covariances proportional to those required by the Goemans-Williamson algorithm. The spiking threshold of the LIF neurons implements a rounding and sampling operation that we map to graph cuts. Neurons that spike together on a given timestep map to vertices on one side of the cut, and neurons that are silent on a given timestep map to vertices on the other side of the cut.

### B. LIF-Trevisan

The second neural circuit (Figure 2) implements the simple spectral modification of Trevisan's algorithm [21], [29] and we refer to it as either LIF-Trevisan or LIF-TR. Like the LIF-GW circuit, the first stage consists of a population of LIF neurons, one for each vertex in the graph, driven by a pool of random devices. Next, the output of the LIF population is fed onto a single LIF neuron. The output of this Stage-2 neuron is discarded; what matters is the weight vector $w$ linking the Stage-1 LIF population to Stage-2. This weight vector is controlled by Oja's anti-Hebbian plasticity rule. This forces the weight vector $w$ to converge onto the minimum eigenvector of the LIF covariance matrix.

The LIF covariance matrix is determined by the connection weights between the random devices and the LIF population. These are set proportional to the Trevisan matrix, which is the sum $I + D^{-1/2}AD^{-1/2}$ of the identity plus the normalized adjacency matrix of the graph. In this way, the LIF-Trevisan circuit does not require solving an SDP offline.

### V. RESULTS

We simulated these circuits and quantified their ability to generate graph cuts. Following [21] we evaluated the circuits on Erdős-Rényi random graphs with a number of vertices $n$ in $\{50, 100, 200, 350, 500\}$ and a connection probability $p$ in $\{0.1, 0.25, 0.5, 0.75\}$. We generated 10 distinct random graphs per $(n, p)$ combination, yielding 200 total graphs. We generated $2^{20}$ graph cuts per circuit, per graph. We compared the circuit-generated cut weights to cut weights generated by a generic SDP solver (PyManOpt [28]) and cut weights generated by a purely random assignment of vertices to sides of a cut. As described, circuits were driven by a simulated pool of random devices. Each device was assumed to have two states, and have a probability of 0.5 of being in any given state at each time step.

Figure 3 shows that, as expected, the LIF-GW circuit matches the performance of the generic solver.
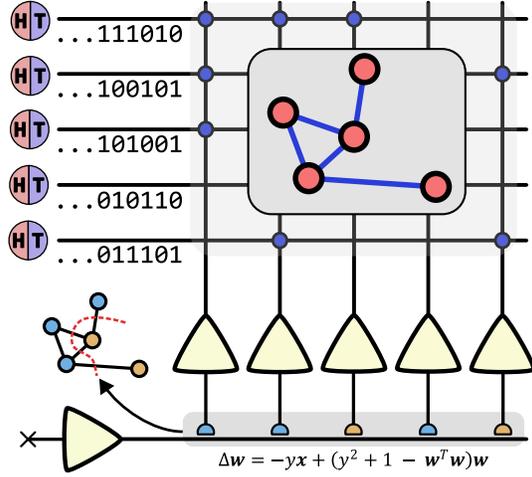


Figure 2. LIF-Trevisan circuit implementing a stochastic approximation to MAXCUT by combining hardware randomness with anti-Hebbian synaptic plasticity. The connection weights between the random device pool (left) and the LIF neurons are set proportional to the adjacency matrix of the graph. The activity of the LIF neurons drives synaptic plasticity on the weight vector onto an output neuron. The solution is sampled by thresholding this weight vector by sign: excitatory, positive weights correspond to one side of the cut and inhibitory, negative weights correspond to the other side. The output of the output neuron is ignored.

This also validates the proposed circuit motif using LIF neurons to translate hardware randomness into Gaussian processes with desired covariances. The LIF-Trevisan circuit shows performance that increases over time, approaching the performance of the solver, due to the on-line learning of the solution through Oja's rule. In all cases, the LIF-Trevisan circuit eventually outperforms the random algorithm. The trajectory of the LIF-Trevisan circuit's performance suggests that the rate of convergence of the plasticity depends on the graph parameters, but there is no indication that the performance of this circuit is saturated within the number of samples considered here.

We next evaluated our circuits on empirical graphs taken from the Network Repository [26]. We picked the same graphs tested in [21]. Figure 4 shows the performance of our circuits compared to SDP-solver and random cuts. Consistent with out results on Erdős-Rényi graphs, we found that the LIF-GW sampling circuit matched the performance of the software solver. We found that the LIF-Trevisan circuit was able to outperform randomly-generated cuts and, occasionally, exactly match the solver-generated cuts, after evolving the circuit with synaptic plasticity for sufficiently many samples. This is consistent with the results of [21], who found that in some cases the simplified approximation algorithms to MAXCUT matched or outperformed cuts generated by the full Goemans-Williamson algorithm on empirical graphs, even though the simplified algorithms have worse approximation guarantees. The maximum cut values
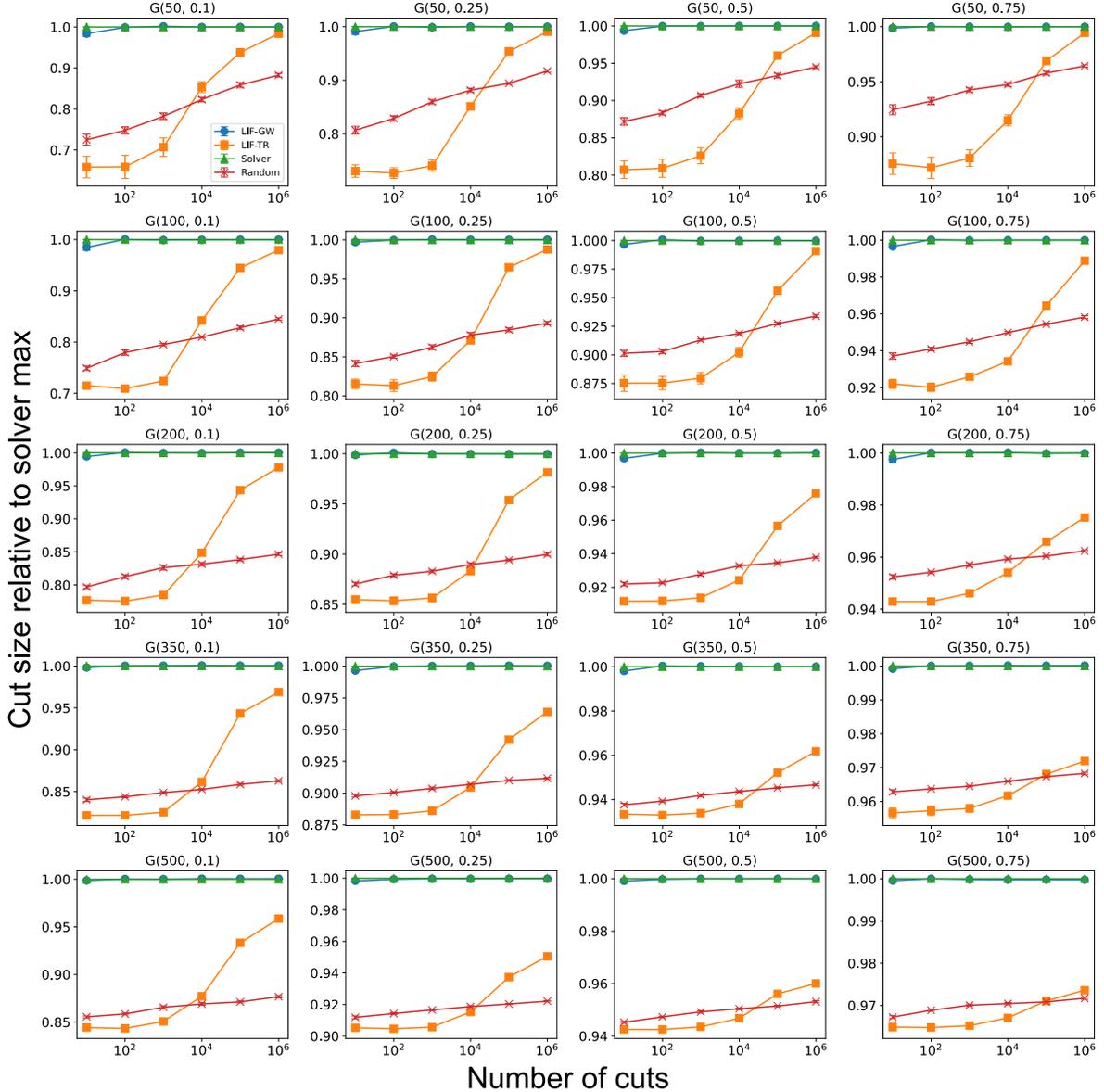
Figure 3. Maximum cut weight relative to software Goemans-Williamson solver (green triangle curve) as a function of the number of samples for Erdős-Rényi random graphs. Rows correspond to fixed numbers of vertices $n$ and columns correspond to fixed connection probabilities $p$. Panel title gives graph parameters $G(n, p)$. Error bars correspond to standard error of the mean over 10 independently generated graphs from each graph class. Blue circles: LIF-GW circuit; orange squares: LIF-TR circuit; green triangles: software solver; red X's: random graph cuts. Blue and green curves overlap.

for each circuit for each graph are presented in Table I, which are in agreement with the maximum cut sizes found in [21] (rightmost column).

## VI. DISCUSSION

We have presented two neuromorphic circuits that transform the activity of a pool of random devices into useful distributions that solve a computational problem, in this case, MAXCUT. Our approach combines insights from theoretical computer science, neuroscience, and materials science to show that probabilistic neural computing is a viable path to new computational architectures. Our circuits display

competitive performance with traditional software solvers. Consistent with prior work [21], we find that the simple spectral Trevisan algorithm performs in practice nearly as well as the gold-standard Goemans-Williamson algorithm. Our results for the Trevisan circuit suggest that its performance can be expected to improve beyond the $2^{20}$ samples considered here. While in this work we considered only a single example problem, MAXCUT is a special case of a larger class of problems known as constraint satisfaction problems, which include problems like maximum directed cut (MAXDICUT) and maximum 2-satisfiability (MAX2SAT). As such, our circuits
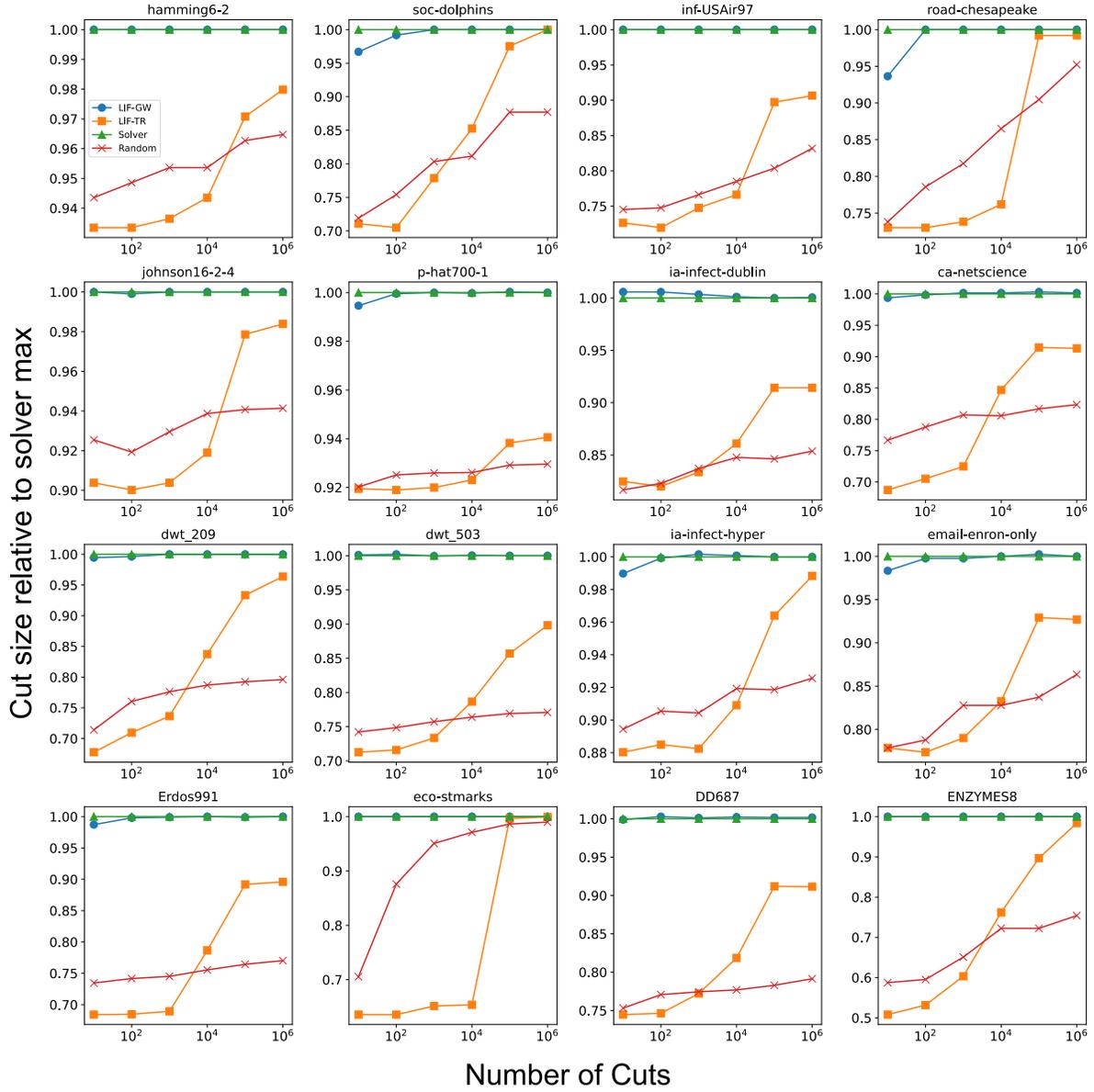
Figure 4. Maximum cut relative to solver as a function of the number of samples for empirical graphs taken from the Network Repository. Each panel represents a single graph, thus there are no error bars. Panel title identifies the graph dataset. Blue circles: LIF-GW circuit; orange squares: LIF-TR circuit; green triangles: software solver; red X's: random graph cuts. Blue and green curves overlap.

may extend to more general probabilistic neural approaches for solving discrete optimization problems. For instance, using results due to Goemans and Williamson [9], our LIF-GW circuit can implement sampling steps for algorithms for MAXDICUT and MAX2SAT that yield approximation ratios of 0.796 and 0.878, respectively.

Our circuits present a trade-off in neuromorphic implementations of combinatorial optimization. The LIF-GW circuit requires fewer random devices and delivers superb solutions rapidly, but requires a substantial commitment of offline resources (i.e. solving a semi-definite program) to initialize. Conversely, the LIF-Trevisan circuit requires as many random

devices as vertices in the graph, and takes many more samples to reach comparable performance. However, this circuit avoids offline computations, which are a significant fraction of the running time of these algorithms [21]. This prior work also used far fewer samples (100) to compare the algorithms than our $2^{20}$. While the number of samples required suggests a disadvantage, at the speed of hardware, the greater number of samples required will likely be a trivial increase in the running time compared to a software implementation. Current hardware implementations of LIF neurons operate with time constants on the order of 1 nanosecond [7], [12]. Using this value as a reference time step for a hardware implementation

| Graph | LIF-GW | LIF-TR | Solver | Random | [21] |
|---|---|---|---|---|---|
| hamming6-2 | 992 | 972 | 992 | 957 | 992 |
| soc-dolphins | 122 | 122 | 122 | 107 | 121 |
| inf-USAir97 | 107 | 97 | 107 | 89 | 107 |
| road-chesapeake | 126 | 125 | 126 | 120 | 125 |
| johnson16-2-4 | 3036 | 2987 | 3036 | 2858 | 3036 |
| p-hat700-1 | 33350 | 31369 | 33351 | 31002 | 33050 |
| ia-infect-dublin | 1751 | 1600 | 1750 | 1494 | 1664 |
| ca-netscience | 635 | 579 | 634 | 522 | 611 |
| dwt-209 | 554 | 534 | 554 | 441 | 540 |
| dwt-503 | 1937 | 1740 | 1937 | 1493 | 1921 |
| ia-infect-hyper | 1277 | 1262 | 1277 | 1182 | 1233 |
| email-enron-only | 425 | 394 | 425 | 367 | 413 |
| Erdos991 | 1027 | 920 | 1027 | 791 | 934 |
| eco-stmarks | 1765 | 1764 | 1765 | 1747 | 1190 |
| DD687 | 1786 | 1625 | 1783 | 1411 | 1680 |
| ENZYMES8 | 126 | 124 | 126 | 95 | 126 |

of these circuits, the circuits could generate millions of samples in the time required for a software simple spectral computation ($\sim$ 10ms), or billions of samples in the time required to solve and sample the Goemans-Williams SDP [21]. The convergence rate of the synaptic plasticity in our LIF-TR circuit depends on both the circuit parameters and graph structure in complicated ways, but this dependence could be formalized or optimized in future work. Furthermore, the requirement for greater numbers of random devices is likely not a limitation, as the current trajectory for implementing the stochastic devices required for these circuits shows promising scaling advantages [22].

Our simulations model random devices as perfectly fair coins generating random, independent bit streams. These assumptions are necessarily approximations to the true behavior of a random device, which may display the statistics of an unfair coin, show internal or external correlations, or display statistics that drift over time. These imperfections might have an impact on the performance of the circuits presented. The key circuit motif in each circuit implements the central limit theorem through the integration of large numbers of random devices. Thus, we expect robustness to deviations of individual devices from the idealized perfect coin as the number of devices grows. While there is a growing realization that stochastic devices can provide robust random bit streams [25], there currently are few standards for what makes a good true random number generator for randomized algorithms. For this reason, the circuits described here provide a much needed benchmark for device physicists to incorporate physically-detailed device models to assess the impact of device variability.

Neuromorphic computing is having a growing impact on graph algorithms [1]. Previous work has found neuromorphic solutions to graph problems such as max flow [17], cycle detection [17], shortest paths [2], [3], [16], and spanning trees [16]. This prior work has exploited connections between the graph structure of neural networks and the corresponding graph problems. In contrast, our work uses the statistical behavior of neural circuits and learning through synaptic plasticity to solve a new class of graph problems. Incorporating learning into neuromorphic circuits to solve specific computational problems is comparatively under-explored, and thus our work expands the neuromorphic acceleration of graph algorithms in a new direction.

Our approach presents a different strategy for incorporating neuroscientific insight into parallel computation. The most successful application of neuroscience principles to date occurs in the field of deep learning, which is based on connectionist principles inherited from early neuroanatomical studies. In contrast, our circuits use the integrative, statistical properties of neurons to achieve different kinds of computations. This was informed by early theoretical developments, like Oja's rule for synaptic plasticity [23], [24]. These developments have been well-known for decades, but have not had as strong an influence on computation. Similar to how backpropagation was known for many years before physical implementations achieved the scale necessary to reveal its utility, we expect that recent advances in stochastic devices and neuromorphic hardware will reveal the utility of probabilistic neural computation.

REFERENCES

[1] James B. Aimone, Prasanna Date, Gabriel A. Fonseca-Guerra, Kathleen E. Hamilton, Kyle Henke, Bill Kay, Garrett T. Kenyon, Shruti R. Kulkarni, Susan M. Mniszewski, Maryam Parsa, Sumedh R. Risbud, Catherine D. Schuman, William Severa, and J. Darby Smith. A review of non-cognitive applications for neuromorphic computing. *Neuromorphic Computing and Engineering*, 2(3):032003, September 2022. Publisher: IOP Publishing.

[2] James B. Aimone, Yang Ho, Ojas Parekh, Cynthia A. Phillips, Ali Pinar, William Severa, and Yipu Wang. Provable Advantages for Graph Algorithms in Spiking Neural Networks. In *Proceedings of the 33rd ACM Symposium on Parallelism in Algorithms and Architectures*, SPAA '21, pages 35–47, New York, NY, USA, July 2021. Association for Computing Machinery.

[3] James B. Aimone, Ojas Parekh, Cynthia A. Phillips, Ali Pinar, William Severa, and Helen Xu. Dynamic Programming with Spiking Neural Computing. In *Proceedings of the International Conference on Neuromorphic Systems*, ICONS '19, pages 1–9, New York, NY, USA, July 2019. Association for Computing Machinery.

[4] Joao Basso, Edward Farhi, Kunal Marwaha, Benjamin Villalonga, and Leo Zhou. The Quantum Approximate Optimization Algorithm at High Depth for MaxCut on Large-Girth Regular Graphs and the Sherrington-Kirkpatrick Model. In François Le Gall and Tomoyuki Morimae, editors, *17th Conference on the Theory of Quantum Computation, Communication and Cryptography (TQC 2022)*, volume 232 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 7:1–7:21, Dagstuhl, Germany, 2022. Schloss Dagstuhl – Leibniz-Zentrum für Informatik.

[5] Ramón Bernardo-Gavito, Ibrahim Ethem Bagci, Jonathan Roberts, James Sexton, Benjamin Astbury, Hamzah Shokeir, Thomas McGrath, Yasir J. Noori, Christopher S. Woodhead, Mohamed Missous, Utz Roedig, and Robert J. Young. Extracting random numbers from quantum tunnelling through a single diode. *Scientific Reports*, 7(1):17879, December 2017.

[6] Dimitris Bertsimas and Yinyu Ye. *Semidefinite Relaxations, Multivariate Normal Distributions, and Order Statistics*, pages 1473–1491. Springer US, Boston, MA, 1998.

[7] Wesley H. Brigner, Naimul Hassan, Lucian Jiang-Wei, Xuan Hu, Diptish Saha, Christopher H. Bennett, Matthew J. Marinella, Jean Anne C. Incorvia, Felipe Garcia-Sanchez, and Joseph S. Friedman. Shape-Based Magnetic Domain Wall Drift for an Artificial Spintronic Leaky Integrate-and-Fire Neuron. *IEEE Transactions on Electron Devices*, 66(11):4970–4975, November 2019. Conference Name: IEEE Transactions on Electron Devices.

[8] Chase Cook, Hengyang Zhao, Takashi Sato, Masayuki Hiromoto, and Sheldon X. D. Tan. GPU-based Ising computing for solving max-cut combinatorial optimization problems. *Integration*, 69:335–344, November 2019.

[9] Michel X. Goemans and David P. Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *J. ACM*, 42(6):1115–1145, 1995.

[10] Hidenori Gyoten, Masayuki Hiromoto, and Takashi Sato. Area Efficient Annealing Processor for Ising Model without Random Number Generator. *IEICE TRANSACTIONS on Information and Systems*, E101-D(2):314–323, February 2018. Publisher: The Institute of Electronics, Information and Communication Engineers.

[11] Hidenori Gyoten, Masayuki Hiromoto, and Takashi Sato. Enhancing the Solution Quality of Hardware Ising-Model Solver via Parallel Tempering. In *2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, pages 1–8, San Diego, CA, USA, November 2018. IEEE Press.

[12] Naimul Hassan, Xuan Hu, Lucian Jiang-Wei, Wesley H. Brigner, Otitoaleke G. Akinola, Felipe Garcia-Sanchez, Massimo Pasquale, Christopher H. Bennett, Jean Anne C. Incorvia, and Joseph S. Friedman. Magnetic domain wall neuron with lateral inhibition. *Journal of Applied Physics*, 124(15):152127, October 2018. Publisher: American Institute of Physics.

[13] D. O. Hebb. *The organization of behavior; a neuropsychological theory*. The organization of behavior; a neuropsychological theory. Wiley, Oxford, England, 1949. Pages: xix, 335.

[14] Giacomo Indiveri, Bernabe Linares-Barranco, Tara Hamilton, André van Schaik, Ralph Etienne-Cummings, Tobi Delbruck, Shih-Chii Liu, Piotr Dudek, Philipp Häfliger, Sylvie Renaud, Johannes Schemmel, Gert Cauwenberghs, John Arthur, Kai Hynna, Fopefolu Folowosele, Sylvain SAÏGHI, Teresa Serrano-Gotarredona, Jayawan Wijekoon, Yingxue Wang, and Kwabena Boahen. Neuromorphic Silicon Neuron Circuits. *Frontiers in Neuroscience*, 5, 2011.

[15] Richard M. Karp. Reducibility among Combinatorial Problems. In Raymond E. Miller, James W. Thatcher, and Jean D. Bohlinger, editors, *Complexity of Computer Computations: Proceedings of a symposium on the Complexity of Computer Computations, held March 20–22, 1972, at the IBM Thomas J. Watson Research Center, Yorktown Heights, New York, and sponsored by the Office of Naval Research, Mathematics Program, IBM World Trade Corporation, and the IBM Research Mathematical Sciences Department*, pages 85–103. Springer US, Boston, MA, 1972.

[16] Bill Kay, Prasanna Date, and Catherine Schuman. Neuromorphic Graph Algorithms: Extracting Longest Shortest Paths and Minimum Spanning Trees. In *Proceedings of the Neuro-inspired Computational Elements Workshop*, NICE '20, pages 1–6, New York, NY, USA, March 2020. Association for Computing Machinery.

[17] Bill Kay, Catherine Schuman, Jade O'Connor, Prasanna Date, and Thomas Potok. Neuromorphic Graph Algorithms: Cycle Detection, Odd Cycle Detection, and Max Flow. In *International Conference on Neuromorphic Systems 2021*, ICONS 2021, pages 1–7, New York, NY, USA, July 2021. Association for Computing Machinery.

[18] Andrew D. Kent and Daniel C. Worledge. A new spin on magnetic memories. *Nature Nanotechnology*, 10(3):187–191, March 2015. Number: 3 Publisher: Nature Publishing Group.

[19] Subhash Khot, Guy Kindler, Elchanan Mossel, and Ryan O'Donnell. Optimal inapproximability results for MAX-CUT and other 2-variable CSPs? *SIAM Journal on Computing*, 37(1):319–357, 2007.

[20] Christof Koch. *Biophysics of Computation: Information Processing in Single Neurons*. Oxford University Press, 1999.

[21] Renee Mirka and David P. Williamson. An Experimental Evaluation of Semidefinite Programming and Spectral Algorithms for Max Cut. In Christian Schulz and Bora Uçar, editors, *20th International Symposium on Experimental Algorithms (SEA 2022)*, volume 233 of *Leibniz International Proceedings in Informatics (LIPIcs)*, pages 19:1–19:14, Dagstuhl, Germany, 2022. Schloss Dagstuhl – Leibniz-Zentrum für Informatik. ISSN: 1868-8969.

[22] Shashank Misra, Leslie C. Bland, Suma G. Cardwell, Jean Anne C. Incorvia, Conrad D. James, Andrew D. Kent, Catherine D. Schuman, J. Darby Smith, and James B. Aimone. Probabilistic Neural Computing with Stochastic Devices. *Advanced Materials*, In Press, 2022.

[23] Erkki Oja. Simplified neuron model as a principal component analyzer. *Journal of Mathematical Biology*, 15(3):267–273, November 1982.

[24] Erkki Oja. Principal components, minor components, and linear neural networks. *Neural Networks*, 5(6):927–935, November 1992.

[25] Laura Rehm, Corrado Carlo Maria Capriata, Misra Shashank, J. Darby Smith, Mustafa Pinarbasi, B. Gunnar Malm, and Andrew D. Kent. Stochastic magnetic actuated random transducer devices based on perpendicular magnetic tunnel junctions. *arXiv:2209.01480 [cond-mat, physics:physics]*, September 2022. arXiv: 2209.01480.

[26] Ryan A. Rossi and Nesreen K. Ahmed. The network data repository with interactive graph analytics and visualization. In *AAAI*, 2015.

[27] José A. Soto. Improved Analysis of a Max-Cut Algorithm Based on Spectral Partitioning. *SIAM Journal on Discrete Mathematics*, 29(1):259–268, 2015.

[28] J. Townsend, N. Koep, and S. Weichwald. PyManopt: a Python toolbox for optimization on manifolds using automatic differentiation. *Journal of Machine Learning Research*, 17(137):1–5, 2016.

[29] Luca Trevisan. Max Cut and the Smallest Eigenvalue. *SIAM Journal on Computing*, 41(6):1769–1786, January 2012. Publisher: Society for Industrial and Applied Mathematics.

[30] Masanao Yamaoka, Chihiro Yoshimura, Masato Hayashi, Takuya Okuyama, Hidetaka Aoki, and Hiroyuki Mizuno. A 20k-Spin Ising Chip to Solve Combinatorial Optimization Problems With CMOS Annealing. *IEEE Journal of Solid-State Circuits*, 51(1):303–309, January 2016. Conference Name: IEEE Journal of Solid-State Circuits.