

Energy-efficient stochastic computing with superparamagnetic tunnel junctions

Matthew W. Daniels,^{1,2,*} Advait Madhavan,^{1,2} Philippe Talatchian,^{1,2} Alice Mizrahi,^{1,2,3} and Mark D. Stiles^{1,†}

¹*Physical Measurement Laboratory, National Institute of Standards and Technology, Gaithersburg, MD, USA*

²*Institute for Research in Electronics and Applied Physics,
University of Maryland, College Park, MD, USA*

³*Unité Mixte de Physique, CNRS, Thales, Univ. Paris-Sud,
Université Paris-Saclay, 91767 Palaiseau, France*

Superparamagnetic tunnel junctions (SMTJs) have emerged as a competitive, realistic nanotechnology to support novel forms of stochastic computation in CMOS-compatible platforms. One of their applications is to generate random bitstreams suitable for use in stochastic computing implementations. We describe a method for digitally programmable bitstream generation based on pre-charge sense amplifiers. This generator is significantly more energy efficient than SMTJ-based bitstream generators that tune probabilities with spin currents and a factor of two more efficient than related CMOS-based implementations. The true randomness of this bitstream generator allows us to use them as the fundamental units of a novel neural network architecture. To take advantage of the potential savings, we codesign the algorithm with the circuit, rather than directly transcribing a classical neural network into hardware. The flexibility of the neural network mathematics allows us to adapt the network to the explicitly energy efficient choices we make at the device level. The result is a convolutional neural network design operating at ≈ 150 nJ per inference with 97 % performance on MNIST—a factor of 1.4 to 7.7 improvement in energy efficiency over comparable proposals in the recent literature.

I. INTRODUCTION

Magnetic tunnel junctions are poised to make significant contributions to new computer chips, most immediately from non-volatile memory applications [1–3], such as magnetic random access memory (MRAM). These devices consist of two magnetic layers separated by a thin tunnelling barrier. The memory values of 0 and 1 are encoded in two different stable configurations of the device (parallel and anti-parallel magnetizations). Values can be read by passing a small current (≈ 10 μ A) through the device, since its resistance depends on its configuration. The device state can be switched by overcoming an energy barrier between the two configurations by passing a higher current through the device. Since MRAM is used as a non-volatile memory, retention requirements demand that the energy barrier be kept high (greater than 40 kT).

However, if the energy barrier is decreased by a factor of ten (to ≈ 4 kT) [4–6], then thermal fluctuations at room temperature cause the device to randomly switch between its stable configurations. In this case, the mean time between thermal switching events is about 55 ns, and the magnetic tunnel junction is said to be in a superparamagnetic state, making such devices superparamagnetic tunnel junctions (SMTJs). In principle, the barrier could be lowered even further for faster switching. The relative times spent in each configuration can be controlled by passing a current through the device creating a spin-transfer torque [7] or by passing a current

through an adjacent heavy metal layer creating a spin-orbit torque [8]. They can currently be fabricated down to a 10 nm length scale [9, 10].

The low energy, truly random behavior, ease of control, and established compatibility with complementary-metal-oxide-semiconductor (CMOS) circuitry has led to the use of SMTJs as the basis for a number of novel computing schemes [11–13]. SMTJs were proposed to implement the concept of probabilistic bits, or p -bits, which were leveraged for applications as Bayesian neural networks [14–16], invertible Boolean logic [17, 18], reservoir computing [19], and Ising network models applied to optimization problems [20, 21]. SMTJs were also proposed as stochastic neural units [22] that can interact with synaptic units, emulated by crossbar arrays of magnetic tunnel junctions that can switch stochastically in the presence of current pulses [23, 24].

Some of these schemes encode information in the switching rates [5] and others encode information in the relative time spent in each state [13, 17]. In these cases, the rates or probabilities are controlled by currents. Though current-controlled methods have been used in many previous device proposals, the ohmic losses they incur can be significant.

Many of these previous applications can be classified as types of probabilistic computing. The specific (and ambiguously named) subfield of probabilistic computing we consider in this paper is called *stochastic computing* [25]. Many of the works mentioned above would not qualify as stochastic computing *per se*. Stochastic computing is concerned with encoding real-valued numbers as the expectation values of random bitstreams; for example, a bitstream such as 01001101000 has four ones and six zeros, thereby encoding the value 0.4, since the probability of seeing a 1 on this wire is 4/10. Stochastic computing

* matthew.daniels@nist.gov

† mark.stiles@nist.gov

is a competitive candidate for energy-efficient application specific architectures [26]. Its robustness to noise [27], high density [27, 28], intrinsic parallelism [26, 27], and latency-precision trade-off [26] make it a promising platform for the implementation of dataflow-based computations in CMOS circuits.

Ideally, stochastic computers operate on long chains of non-repeating, uncorrelated bitstreams that are cheap to produce from an area and energy efficiency standpoint. The conventional way of producing bitstreams is based on a circuit called the linear feedback shift register (LFSR), which comprises a series of flip-flops and simple combinational circuits. LFSRs operate by cycling through all of their internal binary states, each producing a pseudorandom bit, before returning back to the initial state. When LFSRs are used to generate a vast number of pseudorandom bitstreams, those bitstreams are both periodic and cross-correlated. Such correlations are usually undesirable from a computational perspective [26, 29]. Pseudorandom bitstream generation can incur significant overheads in accelerator architectures [28, 30].

Stochastic bitstreams with neither periodicity nor cross-correlation can be generated by SMTJs. The thermal nature of their switching behaviour makes bitstreams generated by SMTJ circuits aperiodic and truly random [31–33]. We replace LFSR-based stochastic sources with arrays of SMTJs that use energy-efficient readout circuitry and programmable logic to generate truly random bitstreams, demonstrating their application in a convolutional neural network architecture.

Using SMTJs to generate stochastic bitstreams becomes useful only when the statistics of those bitstreams can be controlled. Previous works have focused on using current biasing to realize a steady-state spin-transfer torque on the free layer of the junction. This modifies the effective energy landscape of the device so that one of the configurations is tunably preferred over the other. The use of spin-orbit torques in a similar context has also been explored. In this paper, however, we demonstrate how statistical control can be readily accomplished by traditional circuit design; moreover, we show that such digital control has superior energy efficiency over current control for a large range of reasonable material parameters. We discuss the challenges to be overcome, and the contexts to be used, wherein spintronic probability control may become an efficient option.

To make this demonstration concrete, we develop the circuits and architecture to use SMTJs as the basis for a neural network designed to recognize hand-written digits. This allows us to compare, in an application, digitally programmable SMTJ-sourced stochastic bitstreams against spintronically controlled versions. Considering a full scale stochastic computing application also lets us compare SMTJ-sourced stochastic bitstreams against the performance of digitally generated pseudorandom bitstreams. The interplay between high level design requirements for the neural network and the capabilities of the low level devices (SMTJs) leads to modifications

in the design of both the networks and the circuits that connect with the devices. Such engineering across the computational hierarchy, or *stack*, plays an important role throughout this paper. In traditional computational systems, cross-stack engineering is not necessary because clean abstraction layers have been identified between different levels of design to allow optimization of each layer by itself. Computer programmers do not need to know the details of circuits in order to write code, and electrical engineers do not need to know the details of device physics in order to layout useful circuits. These clean abstraction layers break down when using novel devices or novel architectures.

In our case, for example, implementing the neural network with stochastic bitstreams requires uncorrelated bits that can be generated at low energy, a requirement at the device and circuit level driven by the architecture. At the same time, the use of simple primitives, **AND** and **OR** gates, to implement pieces of the neural network dictates changes to the high level structure of the neural network. Similar engineering across the computational stack is important for the many bio-inspired or other alternative computing approaches that aim to take advantage of materials and devices wherein the native dynamics manifest the behavior of neural processes. Because the algorithmic operation depends on the physics (but not vice versa), bottom-up approaches like ours, which choose interesting or energy-efficient physical systems at the foundational level, demand that we pay attention to whether and how the high-level computation can effectively utilize the physics.

In Sec. II, we propose a circuit based on a precharge sense amplifier (PCSA) to read the states of SMTJs to generate a stochastic bitstream. Including a set-reset latch with the PCSA fixes the output to a form useful for stochastic computing. This modified PCSA avoids controlling the state of the SMTJ by write currents; these currents are much higher than the currents needed to read the state, and give rise to ohmic losses that can dominate the energy consumption of the device. Since the expected value of the bitstream is fixed in the absence of tunable write currents, multiple such bitstreams must be combined to produce variable expected values. In Sec. III, we show how to use digital logic to combine these low energy SMTJ-based stochastic oscillators into programmable bitstream generators. These circuits operate at lower energies than LFSRs and other approaches based on SMTJs.

In Sec. IV, we design and simulate a deep convolutional neural network based on LeNet5 [34] to demonstrate the effectiveness of this SMTJ-based approach. The true randomness of SMTJs relaxes constraints on design space considerations for stochastic circuits, allowing us to take advantage of uncommon stochastic computing ideas. We choose an architecture that uses logical **OR** gates as neurons, minimizing area and power expenditure compared to state-machine based approaches. The **OR** gate simultaneously provides both the summation and nonlinear

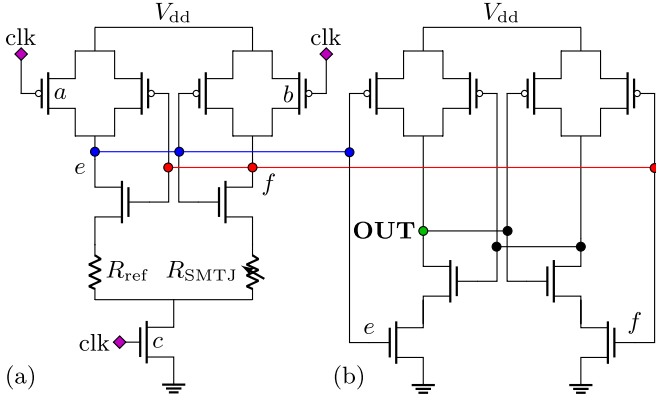


Figure 1. SMTJ readout. (a) The PSCA circuit for reading an SMTJ state [35]. (b) The SC-PSCA includes an set/reset (SR) latch on top of the PSCA from (a), to prevent the pre-charging of nodes e and f from affecting the output duty cycle. The left and right (e and f) branches of the latch are directly wired to the e and f nodes of the PSCA. A more pedagogical version of this circuit is found in Fig. 10.

activation function of a neuron. The use of these devices saves enough energy to justify the modifications needed in the high level architecture. We train the stochastic neural network by backpropagation on an analytic approximation to the network. In Sec. V, we give the results for the accuracy and energy efficiency of this approach based on simulations of this network architecture.

II. PRECHARGE SENSE AMPLIFIER READOUT OF SUPERPARAMAGNETIC TUNNEL JUNCTIONS

Most previous uses of SMTJs use current biasing to vary the duty cycle of the generated bitstream, which can lead to sizeable leakage currents (typically on the order of microwatts). In large scale architectures that require many SMTJs, the resulting ohmic losses can dominate the energy consumption of the computation. An alternate approach using a pre-charge sense amplifier (PSCA) has been proposed in the literature [35]. In this approach, the state is read by a minimal read current pulse and is not controlled by a larger write current.

Others have used the PSCA method successfully to generate random bits [5, 32]. However, these applications only need a single random bit to be produced at a time. In stochastic computing, we need a continuous stream of random bits to be made available at the hardware level. Continuous production happens naturally in current-biased systems where the SMTJ simply sits in a voltage divider, but in the PSCA artifacts of the digital circuitry interfere with extraction of a random bitstream.

Fig. 1(a) shows the PSCA described by Ref. [35]. It works in two cycles. When the clock signal [Fig. 2(d)] is at a low voltage, the transistor at c is turned off but the transistors a and b are on. As there is no path to

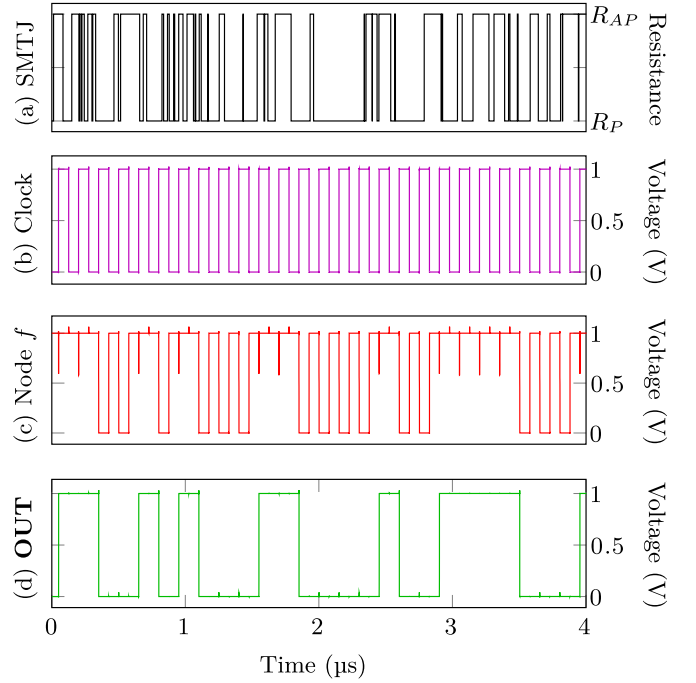


Figure 2. Circuit simulations of the SC-PSCA circuit. (a) Time-series resistance of the SMTJ due to thermal fluctuations. (b) Sampling clock (clk from Fig. 1). (c) Voltage at node f from Fig. 1(a), the output of the standard PSCA circuit. (d) Voltage at node **OUT** from Fig. 1(b), the output of the SC-PSCA.

ground, all wires in the circuit are brought to V_{dd} . When the clock signal goes high, a path to ground is opened at c , and all paths to V_{dd} are closed. Because of the small capacitances C in the transistors, there is a finite discharging time $(R_{sd} + R_{ref})C$ in which charge drains from node e to ground, and a finite time $(R_{sd} + R_{SMTJ})C$ in which charge drains from f to ground, where R_{ref} is the resistance of the reference resistor, R_{SMTJ} is the state-dependent resistance of the SMTJ, and R_{sd} is the source-drain resistance of the transistors between e or f and c . The horizontal red and blue wires provide a nonlinear interaction between these two discharging processes such that the lower resistance channel will connect to ground and the higher resistance channel to V_{dd} after the system comes to equilibrium (Fig. 10). Note that only a small amount of charge $\propto CV_{dd}$ ultimately flows through the system, so ohmic losses are very small; transistor capacitances are typically on the order of 10 aF to 100 aF. Appendix A gives a more detailed discussion of the operation of the PSCA and our proposed modification discussed immediately below.

The problem with the above process is the pre-charge phase, when the clock signal is low and the whole circuit is brought to V_{dd} . The state of the system in that phase does not represent the last measured state of the SMTJ; it is simply preparing to perform the next measurement. This can be seen in Fig. 2(c). Around $t \approx 2.25 \mu s$, for instance, we can see that the SMTJ is in a low resistance

state [Fig. 2(a)]. The output of the PCSA in Fig. 2(c) nevertheless goes to a high voltage repeatedly in this timeframe. This anomalous comb structure on top of the actual SMTJ states is an artifact of the pre-charge phase.

To address this, we attach a circuit called a set/reset (SR) latch to the PCSA design, Fig. 1(b). Whenever e and f are different, the latch copies the values of e and f into its internal wires, so that **OUT** is set to f . When e and f are both brought to V_{dd} during the pre-charge phase, the internal state of the latch is left unchanged. The behavior is explained in detail in Fig. 10 in Appendix A. The simulated state of **OUT** is shown in Fig. 2(d). With the anomalous comb structure removed, this voltage signal becomes suitable for stochastic computing applications.

Fig. 2 shows simulation results of the SC-PCSA obtained using a commercial software package and a 22 nm predictive technology model [36, 37]. The simulations include modeled parasitic contributions from the transistors but not the interconnects, which would depend on layout. The interconnect capacitances would play a role at high speeds, but not the clock periods we consider. In terms of energy, the contributions from the interconnect capacitances are negligible compared to those from the transistor capacitances.

The SMTJ model was implemented in Verilog-A as described in Ref. [38] using parallel and antiparallel resistances of 1.5 k Ω and 4.5 k Ω respectively. The dynamics of the SMTJ are described by the probability to switch in a given time interval, δt of $P_{\text{switch}} = 1 - \exp(-\delta t/\tau)$, where τ is the mean dwell time in that state. It is given by $\tau = \tau_0 \exp(\Delta/kT)$, where Δ is the energy barrier out of the current state of the SMTJ, T is the temperature, k is the Boltzmann constant and τ_0 is a characteristic time scale. In general, the net energy barrier depends on the applied field and voltage, which we do not apply in this paper. These simulations use $\tau_0 = 10^{-9}$ s and $\Delta/kT = 4$.

It is not obvious that SMTJ output sampled at a given time is decorrelated from the output sampled on the same device a single clock cycle later. If the SMTJ is read too frequently, the state of the device has no time to change, and the generated bitstream will have strong autocorrelation. This autocorrelation can be suppressed by increasing the sampling interval. In Fig. 3, we see from numerical simulation that the autocorrelation between two time-adjacent samplings vanishes once the clock interval is chosen to be more than twice the mean dwell time τ of the SMTJ. In other words, the clock cycle should be chosen so that the expected “period” 2τ of a $P \rightarrow AP \rightarrow P$ cycle fits within a single clock cycle. When the SMTJ switches more slowly than the clock, the autocorrelation increases linearly with τ .

Establishing an appropriately large clock cycle is crucial. In many stochastic computing applications, autocorrelation can introduce not only energy inefficiency but functional incorrectness. In the application we present in Sec. IV, for instance, we will purposefully delay

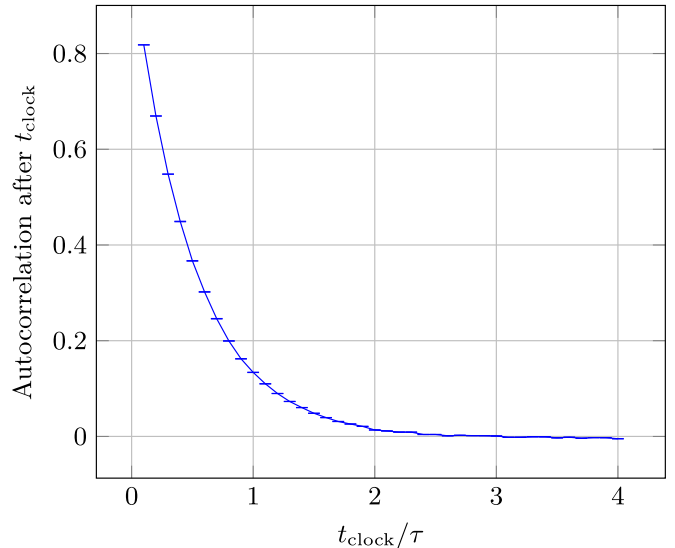


Figure 3. Autocorrelation under a lag of one clock cycle as a function of clock cycle time t_{clock} compared to the SMTJ mean dwell time τ . Each point gives the mean autocorrelation time over 10^3 trials, with each trial sampling the SMTJ for 10^3 clock cycles. The vertical error bars give 95 % confidence intervals on the mean.

time signals from each other in order to suppress cross-correlations among them. This only works if the initial signals themselves have vanishing autocorrelation.

If the continued development of SMTJ technology could engineer a sufficiently small autocorrelation time for the thermally-induced magnetic dynamics—that is, a speed comparable to the switching speed of the CMOS gates—then the integrated ohmic power loss will be similar to losses in the CMOS itself, and it may then become reasonable to place a *static* read current across the devices, as proposed in Refs. [12, 17, 19]. This would provide a continuous-time random telegraph signal, opening up the possibility of running asynchronous computations. We compare the energetic performance of such an approach to our proposal in Sec. III, in the context of current SMTJ technology. We explore this limit in detail in Appendix C.

Since the SC-PCSA is charge based, the dwell time of the SMTJ does not change the energy expenditure of the circuit (so long as the mean dwell times are not faster than the equilibration time of the PCSA, about 1 ns). Ideally, then, we would like the SMTJs to fluctuate as quickly as possible; in general, we would also hope for uniformity of dwell times and operation ranges. Measured dwell times range between 1 μ s and 0.1 s depending on the operating regime [6]. The fluctuation rates are highly sensitive to applied fields and currents. The two state fluctuator model used in this paper is not valid for fluctuator frequencies approaching the 10^{-9} s time scale of magnetization reversal [39], but for applications like ours in which a reference resistor is present, Kaiser et al. [40] show that frequency scales can exceed the 1 GHz

regime. The achievable time scale depends on the differences in the resistance, proportional to the tunneling magnetoresistance of the magnetic junction, that can be detected and the size of the current needed to do so.

The resistance and magnetoresistance of the tunnel junction are the same as those developed for memory applications and so should have margins sufficient for effective circuit design [41]. In memory applications, however, the state of device is switched by the spin transfer torque (STT) that is generated by the current passing through the device. Substantial work has been done to make that switching current as low as possible [42]. In the present application, by contrast, we want the current response to be as weak as possible, so that the read current has a minimal effect on switching rates. Unfortunately, the current needed to influence the switching behavior apparently [9] scales to smaller values as devices become smaller. Increasing the this current requires additional research; recent work suggests that devices with easy-plane anisotropy may be preferable, in this sense, compared to perpendicularly magnetized devices [43].

There has been considerable work in the recent literature to use MTJs or SMTJs as the fundamental units for truly random number generators [31–33, 44]. For general purpose applications, random number generators are required to pass certain tests of randomness; the NIST statistical test suite [45] is usually taken as a standard benchmark for validating good randomness in that sense. An important criterion required by that test suite is that the mean value of a bitstream has a long time average of 50 %. To date, all MTJ-based solutions that pass those tests require XOR'ing eight devices together to eliminate bias. Given SMTJs with similar enough dwell times in each state as we assume here, we expect that we could combine SC-PCSAs in a similar manner and pass the statistical tests needed for random number generation.

In this work, we do not subject our SC-PCSA to these randomness tests, as they are not necessary for our application space. Stochastic computing aims to encode values in the expected value of random bitstreams. It is traditionally implemented with low bit-resolution pseudorandom number generators that have poor randomness properties; some stochastic computing has even been done with entirely deterministic bitstreams [46, 47]. The important properties needed for stochastic computing are simply small enough cross-correlation between different devices and small enough autocorrelation after some time has passed. The correct metric for the success of our generator is not that it can produce unbiased random bits for cryptographic or scientific applications, but that it can be used to carry out stochastic computing calculations. We show that it can in a complex neural network architecture in Sec. IV.

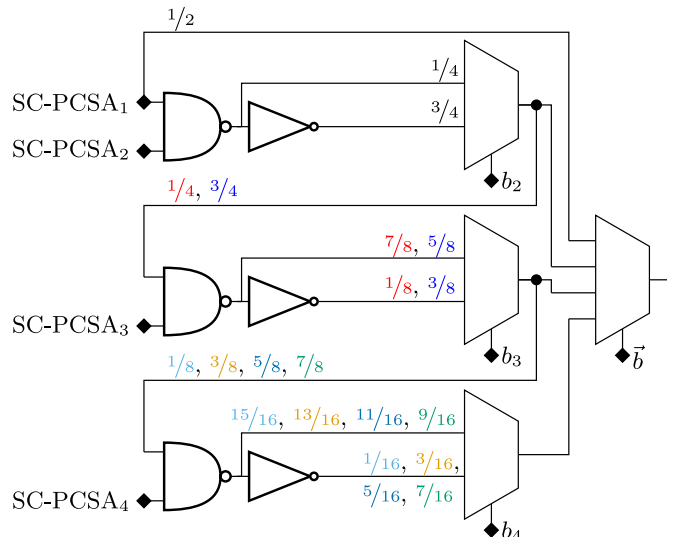


Figure 4. Bitstream generator with four programmable bits. Each row is a recursive subdivision unit. Assuming input probabilities of $1/2$ from the SC-PCSAs, the topmost 2-input multiplexer outputs $1/4$ or $3/4$ depending on b_2 ; the middle multiplexer $1/8$, $3/8$, $5/8$, and $7/8$, depending on b_2 and b_3 ; and so on. SRAM control-and-decode circuits, which set the bits \vec{b} appropriately for the desired output probability, are hidden for clarity.

III. SMTJ PROGRAMMABLE BITSTREAM GENERATOR

In order to generate arbitrary bitstreams with n -bit precision, we develop a composite cell that can be programmed based on values stored in static random access memory (SRAM) cells. If the programmed value is $1/2$, we need only tap the output of one of the SC-PCSA cells discussed in Sec. II. Otherwise, we perform a sort of binary search on the unit interval. Given a stochastic signal carrying the probability value $k/2^{n-1}$, it is synchronously fed into a NAND gate together with a $1/2$ probability signal from a new SC-PCSA cell to produce a signal with probability $1 - k/2^n$. A simple multiplexer (MUX) can then optionally route the signal through an inverter, allowing us to access values of both $1 - k/2^n$ and $k/2^n$.

Iterating this recursive subdivision operation allows us to access all multiples of $1/2^n$ using n SC-PCSA cells. Values for zero and unity can be trivially implemented on the end of the circuit by replacing the stochastic signal with a constant voltage. In our simulations of a neural network later in this paper, we restrict ourselves to the 4-bit case, as illustrated in Fig. 4, allowing our synaptic weights to express integer multiples of $1/16$.

Though strict tests of randomness are not necessary for our purposes, careful attention to correlation is. Statistical correlations can arise in many ways in a stochastic computing circuit. We distinguish between two types: *graph correlation*, and *source correlation*. To understand the former, suppose a circuit designer makes a simpli-

fying assumption that the inputs to a particular circuit node represent statistically independent stochastic processes. This may not be true, depending on whether the computations producing those inputs overlapped at some point earlier in the circuit. If they are not statistically independent, the circuit output statistics will differ from the value expected by the designer. We call the correlations leading to such errors graph correlations, as they arise due to reconvergent fanout in the computational graph of the circuit. One method for alleviating these correlations is called isolation, and has recently been discussed in Ref. [48]. We make use of these and discuss their implementation and impact in Sec. IV.

The programmable bitstream generator described above addresses a more insidious form of correlation that can arise when the original bitstreams are not random. In typical implementations of stochastic computing, the LFSRs and other commonly used sources of pseudorandom bitstreams are periodic. Overusing the same LFSR design at multiple instances in a circuit can lead to correlation-induced errors that are difficult to predict [49]. We call these source correlations; they are induced by inconsistent assumptions about bitstream generators on the boundary of the computational graph. One of our motivations for SMTJ-based stochastic computing is to solve the problem of source correlations. This is especially important in neural networks, where hundreds or thousands of inputs can all fan-in to the same circuit node. To avoid source correlations in this case, pseudorandom circuit designs become prohibitively energy-hungry, but our SMTJ-based approach scales to any degree of fan-in with no source correlation present.

Energy numbers for our approach at various bit-precisions and power supply voltages are shown in Fig. 5. Though lower supply voltages should generally lead to lower energies, we see that for $n = 8$ the energy per cycle is actually higher for the lower supply voltage. In order to get high reliability operation of the SC-PCSA at low supply voltage, the drive strength of some transistors needs to be increased. This results in the 0.8 V SC-PCSA requiring more energy than its 1 V counterpart. As the number of SC-PCSA in the readout circuit increases at high bit-precision, this anomalous energy cost at low supply voltage can overwhelm the nominal savings in the rest of the logic tree.

Fig. 5 also compares SC-PCSA energy to the energy consumption of an LFSR with the same bit precision. For a fair comparison, we designed the (Galois type) LFSRs according to standard maximum sequence length designs, and using the same technology node and predictive technology models as for the SC-PCSA simulations. We found that the SC-PCSA performed about twice as efficiently across a range of bit precisions. This result is predictable from the transistor count of the two designs; the SC-PCSA has 15 transistors, whereas a clocked, NAND-gate based digital flip flop (the fundamental unit of the LFSR) has 20. The circuitry of comparator in the LFSR-based is also slightly larger than the AND-NOT-MUX rows

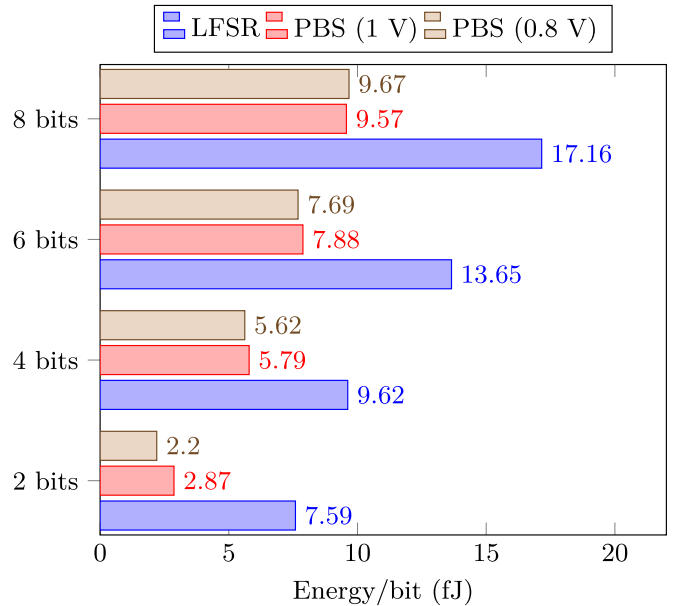


Figure 5. Energy efficiency comparison between our programmable bitstream generator (PBS) at supply voltages of 1 V and 0.8 V and a traditional LFSR with binary comparator. Plotted horizontally is the energy needed to produce a single new element of the output stochastic bitstream. The N -bit PBS uses N SMTJs where the N -bit LFSR uses an N -bit register, ensuring a fair equiprecision comparison.

of Fig. 4. The LFSR-based scheme additionally contains some small number of XOR gates, though the number of XORs is an algebraic property that does not necessarily scale with the number of bits, N .

We note that several proposals address the cost of LFSR-based pseudorandom number generator by extracting combinatorial subsets of the bits in a large LFSR [50, 51]. Though these shared-LFSR methods amortize the energy cost by sharing it over many pseudorandom bits, the resulting correlation between these bitstreams is 1.5 to 2 times higher than the isolated LFSR case [51]. In correlation-sensitive applications, this is clearly disadvantageous. In applications that have been engineered to be correlation insensitive, one could imagine applying the same shared generator techniques to our programmable bitstream generator; the same energy/correlation trade-off should apply. We leave the details of such a device to future research.

Several related configurations for using magnetic tunnel junctions have been proposed that could be used for generating stochastic bitstreams. Most fall into two broad categories, those based on nominally stable magnetic tunnel junctions that are brought into an unstable state by a current [31, 44, 52] and those based on SMTJs that are current biased to control the expected value of the bitstream [14–17, 19–21]. The write currents in both of these approaches lead to significant ohmic losses, which are disadvantageous for this stochastic computing. We note however that in many cases, these configurations

have been proposed for applications for which the ohmic losses may not be quite as important.

We estimate the ohmic losses for pulsed destabilizing of stable MTJs in Appendix B. The circuits reported in Ref. [44] are based on MTJs with average resistance $R = 1000 \Omega$, and use an average write voltage of 217.5 mV for 8.75 ns. The ohmic losses are then 414 fJ per bit. Ref. [52] uses similar techniques on MTJs specifically tailored to the generation of stochastic computing bitstreams (the same task we consider here); they report average costs of 526 fJ per bit. Both estimates are significantly higher than the 10 fJ estimated for the SC-PCSA-based approach.

To estimate the ohmic losses (see Appendix C for details) for p -bit style current-biased SMTJs, we assume a supply voltage of $V_{dd} = 1$ V and favorable resistances of $R_{AP} = 100$ k Ω and $R_P = 50$ k Ω , and find ohmic losses of about 500 fJ per bit for the 150 ns clock cycle we consider; the clock period is set by the autocorrelation time of the SMTJ. In the SC-PCSA-based approach described above, the energy per bit of 10 fJ for the whole circuit is roughly independent of the clock cycle. Note that, for the current controlled approach, the energy per bit decreases linearly as the clock cycle decreases, decreasing the advantage of the SC-PCSA-based approach.

We cannot expect the manufacturing margins of real SMTJ devices to be ideal. The fabrication of magnetic tunnel junctions has been optimized in commercial fabrication facilities [53, 54] in the nonvolatile regime that is desirable for memory applications. SMTJs have only been studied for larger devices fabricated in laboratory settings [6].

Real devices will have distributions of all of their properties. Three important ones are variations in barrier heights (and hence characteristic dwell times), variations in biasing around the perfect $p = 1/2$ that we have assumed above, and variations in the device resistances. The insensitivity to a distribution of barrier heights is addressed in Fig. 3. Essentially uncorrelated bitstreams are generated as long as the mean dwell times are shorter than the clock cycle and longer than some lower limit set by the response time of the PCSA, on the order of 1 ns for the circuits considered here. Assuming a prefactor of $\tau_0 \approx 1$ ns in the transition time distribution, this approach is insensitive to fluctuations in the barrier height for $0 < \Delta/kT < 5$. Since considerable effort has been required to keep energy barriers large as MTJ devices become smaller [55, 56], it seems that it should not be difficult to fabricate devices with small enough barriers to enable thermally driven transitions. Variations in the barrier height should not pose a difficulty if it is possible to maintain the roughly 10 % variation in barrier heights [57] as devices are scaled down.

Having an expected value of 0.5 for the output of the bitstream generator requires that the energies of the parallel and antiparallel states are the same. Such equality in turn requires that the fringing fields acting on the free layer be close to zero. It is difficult to fabricate devices

with precisely these fields and there will naturally be a distribution of relative energies and hence dwell times. Fig. 6 gives the effect of the variation in the relative dwell times on the output of a four-bit stochastic bitstream generator. We assume that the fluctuations in the relative dwell times in the parallel and antiparallel states are distributed around zero with standard deviations given in the figure (see Appendix D for details). These distributions should be compared to the expected distribution for perfectly balanced bitstreams that are sampled for a finite sample size. For the results presented below based on using bitstreams of length 128, the expected distributions of values found from perfectly balanced bitstreams are similar to the distributions found for 10 % variations in the distribution of the relative dwell times.

Finally, we address whether the distributions of resistances would adversely affect the performance of our bitstream generator. We have avoided the use of spintronic control over the random telegraph noise, so resistance variations do not directly affect the probability of our bitstreams. The main failure mode would arise from resistance variations, that is, if the reference resistor R_{ref} , selected before devices are grown and tested, were to not fall between the resistances R_P and R_{AP} of the parallel and antiparallel states.

Recent data from fabrication facilities [58, 59] show that the the manufacturing margins should be tight enough to accomplish this. The 3σ variability in the parallel state resistance is under 10 %. The 3σ variability in the tunneling magnetoresistance is only 2 % to 3 %. If we target a 1.5 k Ω resistance for R_P , and a 4.5 k Ω resistance for R_{AP} , then a reference resistance of 3 k Ω is about 15σ away from both resistances, providing an adequate manufacturing margin even if the magnetoresistance were much smaller than the assumed 200 %. In addition, neural network applications, which we have considered in some detail in the paper, can be robust to single-device failures, if the failure of particular devices can be detected (which it is easy to imagine doing here).

The data we cite above is for high-barrier magnetic tunnel junctions developed for magnetic random access memory (MRAM) but manufacturers are pushing toward the superparamagnetic regime; recently announced [60] work on embedded MRAM devices works with retention times on the order of seconds to milliseconds. This is not too far from the SMTJ speeds fabricated in academic labs, and the reported fabrication distributions continue to be satisfactory in this lower barrier regime. Based on these reports, our proposed approach should be resistant to device fluctuations within the range of fluctuations that can be expected from a dedicated fabrication process.

IV. APPLICATION TO NEURAL NETWORKS

Artificial neural networks, initially inspired by the structure of the brain [61], have become one of the most

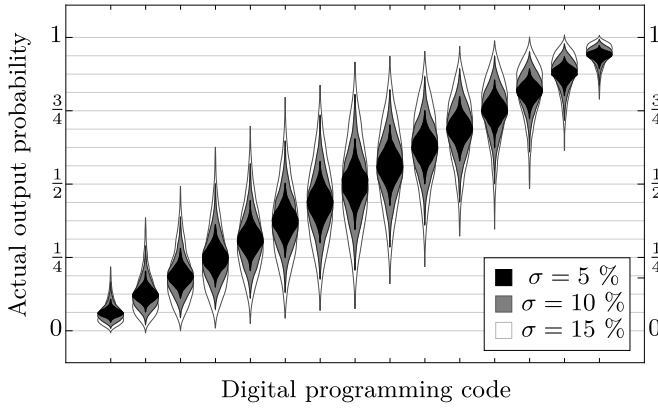


Figure 6. Output probability distributions for the four bit generator induced by imprecision in SMTJ fabrication. The three different standard deviations correspond to the three distributions in Fig. 11. The width of a bubble at a particular vertical coordinate gives the probability density that the generator outputs that probability in its bitstream. Note that the three, two, and one-bit versions can be extracted by simply restricting this plot to outputs with mean values that are multiples of $1/8$, $1/4$, and $1/2$, respectively.

powerful classes of algorithms in machine learning. These neural networks are composed of two fundamental units: neurons, which sum inputs and apply a nonlinear activation function to the resulting sum; and synapses, which multiply the output of one neuron by a real-valued weight and pass the result to a downstream neuron. Traditionally, the numerical values that propagate through neural networks correspond to rates at which neural spikes propagate through an actual biological brain.

The origins of stochastic computing lie in the observation that time series data of stochastic spike trains in the brain could be modeled by stochastic jumps from ground to V_{dd} in a logic circuit [62–64]. It is no surprise, then, that neural network structures have been implemented successfully and energy efficiently in recent stochastic computing work [65–68]. Rather than carrying out high level arithmetic and logic operations to “theoretically predict” a neural network’s output, stochastic computing implements neuromorphic models of the network in CMOS circuitry. The network operation is “experimentally simulated” by the physics of the circuit, and the results are obtained by monitoring the time series voltages at the network output.

A crucial synergy involved in this scheme is the inherent parallelism of both neural networks and stochastic computing. Whereas the many mathematical operations involved in a neural network layer would need to be computed serially in a traditional computing environment, the physical nature of the stochastic computer means that these operations are all run simultaneously. There is a sizeable body of recent work that uses this principle to build efficient, stochastic-computing-based neural networks. The first stochastic computing implementation of a deep convolutional neural network, a particularly im-

portant neural network topology with broad applications to image processing, was proposed by Ren *et alia* [65] and further optimized in Refs. [66, 69]. These works draw heavily on new ideas in the stochastic computing literature, including massively parallel generation of pseudo-random bitstreams [70], state-machine based nonlinear activation functions [67, 71], and aggressive use of correlation insensitivity [72].

A common approach in the stochastic computing neural network literature is to construct a neuron unit from multiplexer-based addition composed with a state-machine-based nonlinearity. These operations are costly in terms of CMOS transistor counts. We take a simpler approach: our entire neuron is a single logical OR gate.

In our approach, a single multi-input logical OR gate simultaneously, approximately, and inseparably performs both the summation and non-linear activation. Though for small input magnitudes an OR gate performs addition of probabilities $p_1 + p_2 + O(p^2)$, nonlinear corrections become important as the input probabilities increase. Utilizing this property to our advantage, we harness that nonlinearity directly as our neuron’s nonlinear activation function. This approach was originally proposed in [73–76], but until now has been unable to scale beyond small networks due to correlations between LFSRs. Our use of truly random bitstreams sourced from SMTJs opens the possibility of using this very small and efficient neuron in modern-scale neural networks. The synapses in our work are AND gates, which naturally implement multiplication on probabilities, with inputs from the output from the previous neuron and from a fresh stochastic bitstream encoding the weight (in our case, provided by an SMTJ programmable bitstream generator).

The probabilistic response of an OR gate in stochastic computing is analogous to its Boolean response; the output probability is given by

$$P_{\text{OR}}(\mathbf{p}) = 1 - \prod_{j=1}^N (1 - p_j), \quad (1)$$

where \mathbf{p} is the vector of input probabilities. Eq. (1) holds when each input probability p_j is an independent random variable. Note that this condition generally fails in stochastic neural networks that use LFSR-based bitstream generation, especially when the number of inputs to a neuron approaches or exceeds the periodicity of the LFSR design being used.

Unlike a traditional artificial neuron, the OR gate’s action on probabilities cannot be factored into two separate processes of summation and activation. But by examining its limiting cases, we can see that the OR gate does *approximately* perform summation-and-activation functionality: for the large fan-in case, where $N > 10$, the lower bound of the OR gate output approaches $1 - \exp(-\sum_j p_j)$. We plot this bound, and all of the other allowed output values, as a function of the input probability sum in Fig. 7. The functional form of this distribution is similar to the hyperbolic tangent function, a

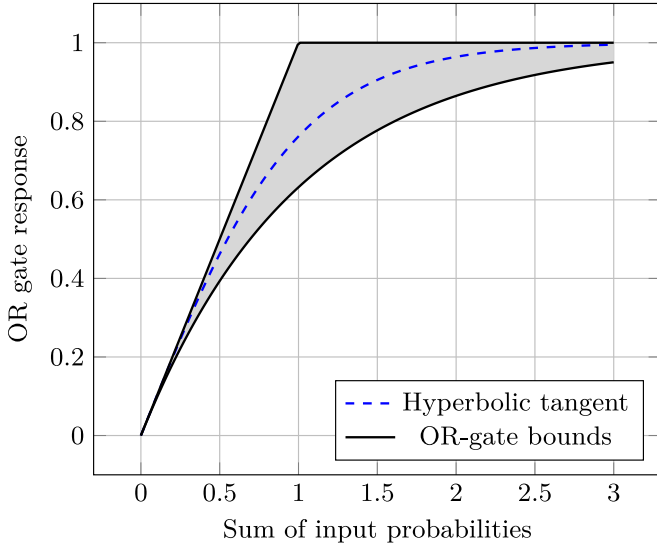


Figure 7. Activation pseudo-function of the OR gate in terms of its bounds. The shaded region indicates possible output values in the large fan-in case. For comparison, the dashed line plots the hyperbolic tangent function, a common neural activation in deep learning applications. Like the hyperbolic tangent, the OR gate operation is asymptotically linear near the origin and saturates exponentially to unity.

widely used activation function in machine learning.

The inputs to a neuron are multiplied by real numbers often called synaptic weights. We represent each weight using a programmable bitstream generator, and each multiplication using an AND gate. The OR-gate neuron, together with AND-gate synapses, form the fundamental unit of our proposed neural network architecture. This is an extremely energy efficient primitive cell, and we have to accept multiple constraints in order to use it. The most striking constraint is that, since all values in our neural network are represented by probabilities between zero and one, our network nominally lacks any form of inhibition (classically implemented with negative numbers) as reflected in the fact that the activation function in Fig. 7 is not defined left of the origin.

Without making changes at the algorithmic level to accommodate the lack of inhibition, the OR-gate neuron would be largely useless. To facilitate the use of this neuron, we replace the standard weight matrices by two sets of weight matrices, one for inhibition, w^I , and one for excitation, w^E , that are merged to provide a form of inhibitory activation for the network. The inhibitory subneuron uses a NOR gate rather than an OR gate so that large inhibitory input suppress the neuron output. The binary output $y \in \{0, 1\}$ of each neuron is given by

$$y = \underbrace{\left[1 - \prod_j (1 - w_j^E x_j) \right]}_{\text{OR}} \times \underbrace{\left[\prod_j (1 - w_j^I x_j) \right]}_{\text{NOR}}, \quad (2)$$

where the outermost multiplication is implemented with

a single AND gate. The use of both excitatory and inhibitory subnetworks introduces a higher dimensional symmetry, which the network can break to accomplish dynamical inhibition (or excitation) and successfully learn correct weights for classification. More details are discussed in Appendix E2.

Since the OR gate neuron is unusual, the network training algorithms must be adapted to it. We derive the backpropagation equations for this neuron and the inhibitory/excitatory subnetworks for use in training, and adapt standard training algorithms to accommodate the hard constraints on allowed weights. The details of this codesign process are elaborated upon in Appendix E. One important note is that Eq. (2) prescribes the correct boolean output at a given clock cycle, but is formally inaccurate as a prescription for the expectation value of y over many cycles. It treats the inputs x to the OR and NOR gates as independent random variables, when in fact they are identical and thus perfectly dependent on each other. Nevertheless, we found empirically that the use of Eq. (2) for network training performed well enough for our purposes. Refining the training process to use a more accurate expression could improve network performance.

To demonstrate the effectiveness of our stochastic neuron and synapses, we implement a stochastic approximation of LeNet5 [34], a convolutional neural network. Formally, LeNet5 is based on floating point arithmetic and specific choices of neurons. We instead use the AND and OR gates for synaptic and neural operations, as outlined above, in a stochastic computing framework. We also make simple changes to the layer structure to account for our restriction to nonnegative probabilities, which are outlined in Appendix E. This benchmark task allows for an easy comparison to existing literature. Ref. [66] treats the standard form of LeNet5 (using stochastic arithmetic) as a fixed boundary condition at the top of the stack, and optimizes a scaffolding of stochastic computing hardware around this software-oriented neural network design. The implementation of LeNet5 we present here is instead based on the kind of cross-stack reasoning we just outlined: we relax the strict definition of LeNet5 as a boundary condition and use the resulting flexibility to accommodate our energy efficient choice of neurons and synapses at lower levels.

V. EVALUATION AND RESULTS

To evaluate our implementation of LeNet5, we train and test our neural network on the MNIST dataset. The MNIST dataset comprises 60000 training images and 10000 test images of 28×28 pixel images. First, we perform offline training on an analytic model based on probabilities rather than stochastic representations of probabilities, that is, on a traditional software model using the unusual topology, constraints, and activation functions of our proposed network. Though the network does train, we find the process to be significantly noisier than train-

Table I. Energy breakdown on worst-case analysis of LeNet5 case study for an SMTJ-driven computation. The rightmost column corresponds to the $N = 128$, $\delta = 16$ network configuration.

| | Energy/fJ Element cycle | Energy/nJ LeNet5 cycle | Energy/nJ inference, at 97 % accuracy |
|-------------------|----------------------------|---------------------------|--|
| $V_{dd} = 1$ V | | | |
| Weights | 5.80 | 0.50 | 112.80 |
| Neurons | 3.11 | 0.04 | 8.09 |
| Isolators | 19.75 | 0.11 | 25.72 |
| Total (nJ) | | 0.65 | 146.62 |
| $V_{dd} = 0.8$ V | | | |
| Weights | 5.62 | 0.49 | 109.39 |
| Neurons | 2.47 | 0.03 | 6.44 |
| Isolators | 15.02 | 0.09 | 19.56 |
| Total (nJ) | | 0.60 | 135.40 |

ing of classical LeNet5 models under the same hyperparameters. We speculate that the noise arises from the unfactorability of the summation-nonlinearity approximated by Eq. 1, as well as the hard constraint that all weights must be between zero and one.

Among successfully trained networks, we find that we frequently achieve around 98 % test set accuracy with the analytic model applied to MNIST, within 1 % of the 98.9 % accuracy achieved by the original LeNet5 authors [34]. We then discretize the trained weights and biases by rounding them to the nearest multiple of $1/16$. These are loaded into a stochastic simulation of the logical architecture, complete with stochastic processes to generate the SMTJ statistics.

To deal with residual graph correlations induced by reconvergent fanout of the signals into and out of each neural network layer, we insert what we call an isolator mask—an array of randomly chosen (but static) pixelwise temporal delays—after each network layer. We discuss in more detail in Appendix E 4 how isolators can be used to achieve decorrelating effects by shifting signals relative to each other in time. Each pixelwise, integer-valued delay is chosen uniformly at random from the interval $[0, \delta]$, where δ is the maximum delay. In our LeNet5 architecture, this leads to a total of $L = 6$ decorrelation layers. The maximally delayed route from input to output, then, has delay length $W = L\delta$, which is the amount of time we must wait before collecting N usable data points at the output.

In the upper black curve on Fig. 8, we plot the performance of a trained network on the test datasets as a function of the maximum delay δ and length of collected bitstream N . Performance improves exponentially with increasing δ , eventually saturating around a maximum delay length of 16—and, therefore, a mean delay length of 8. We also find, in the lower curve, that increasing N provided diminishing returns beyond $N \approx 128$. We use this $(\delta, N) = (16, 128)$ configuration to report the rest of the numbers in our paper. The fluctuations that become

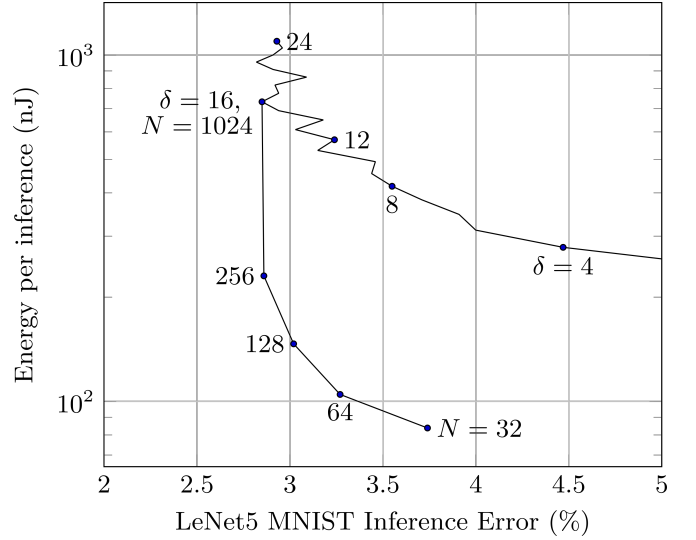


Figure 8. Evaluation of our network performance for various hyperparameters. The upper curve shows variation in δ , the maximum number of isolators decorrelating each of the L hidden layers. This affects the energy by increasing the total warm-up time $L\delta$ needed before network output becomes meaningful. The lower curve tracks N , the number of meaningful data points collected *after* the warm-up period concludes.

evident in the plot at high δ arise from the pseudorandomness of the isolator mask configuration. A systematic approach would search through multiple mask configurations until an appropriate one is found for the given trained network.

To compute the energy efficiency of our proposed architecture, we extract the mean activity factors from weights and neurons in the network over an inference pass on the MNIST test data set. As our simulation is purely logical (that is, it works at the level of clocked bits rather than dynamical voltages), we find it unwieldy to extract the expected activity factors from inside the circuit-level composite gates of the OR gate neuron. Instead, we convert the activity factors from the logical model into the *worst case* activity factors for the realistic architecture.

To produce our worst-case analysis, we assume that *every* switch in the voltage signal causes a maximal number of switches in the AND gate synapse *and* the downstream OR gate neuron. We apply this assumption to the activity factors found in our stochastic model, together with circuit power estimates based on predictive technology models [36, 37]. The results for the $N = 128$, $\delta = 16$ case are listed in Table I. For $V_{dd} = 1$ V at the 22 nm node, CMOS switching events in OR gate neurons and AND gate synapses together account for 8.09 nJ.

Random bitstream generation and isolator-based decorrelation are more expensive than the computational graph. For these estimates we again use worst-case activity factors. In the full LeNet5 architecture, the pre-charge sense amplifier circuits generating network weights and inputs account for about 113 nJ for a full

inference. The programmable isolator buffers are each relatively expensive, but fewer were needed compared to the weights. The total cost for the isolator circuits is around 26 nJ.

We also test the behavior of our architecture at a reduced supply voltage. For voltages below 0.7 V, the PCSA performance begins to degrade, causing errors in the output bitstream. At 0.8 V, the total energy consumed by the entire architecture drops from 147 nJ to 135 nJ. These energy estimates are also displayed in Table I.

One likely source of network inference error that remains after accounting for latency and isolator length is the quantization of weights from the analytic model to discrete values allowed by the bitstream generator. Co-design of hardware with deep neural networks has made the study of low-precision neural networks a topic of considerable recent interest [77]; quantization-aware training methods are now being developed both for pure CMOS systems [78–80] and platforms with nanodevice integration [81, 82]. Applied to systems like ours, these may offer the possibility of improved inference performance at fewer and fewer bits, improving energy consumption and network latency.

With that said, the robustness of our stochastic model after rounding the weights of our analytic model expresses a robustness to error characteristic of neural networks. We speculate that, since network performance was robust with respect to rounding weights to multiples of $1/16$, it should be equally robust had the weights been rounded to values slightly different than the multiples of $1/16$. If that is indeed the case, then systems like ours would enjoy a degree of robustness against variation of the type shown in Fig. 6.

Although some designs for modern neural networks can work at four bit (or sometimes lower) precision, it does not necessarily follow that one could successfully use an $N = 16$ stochastic computing network. Stochastic computers have two types of precision that are intrinsically linked: representational precision, and sample precision. Collecting N time steps means that N different values are representable on the output, but it also means that our certainty (standard deviation) that the measured value matches the true value of the output distribution scales as $N^{-1/2}$. Therefore the expected length of a computation itself influences the meaningful representational precision of the inputs; we found that our 4-bit programmable bitstream generator gave sufficient resolution at $N = 128$. The link between these two types of precision is subtle [83] but will be increasingly important in determining when stochastic computing is or is not energy efficient for different target applications.

In Fig. 9, we compare with the current state of the art in stochastic computing research. The works we cite all treat LeNet5 on MNIST. However, these works were also simulated at the 45 nm, rather than 25 nm, technology node. As is standard practice, we scale our reported energies by the corresponding scaling in transistor size (and

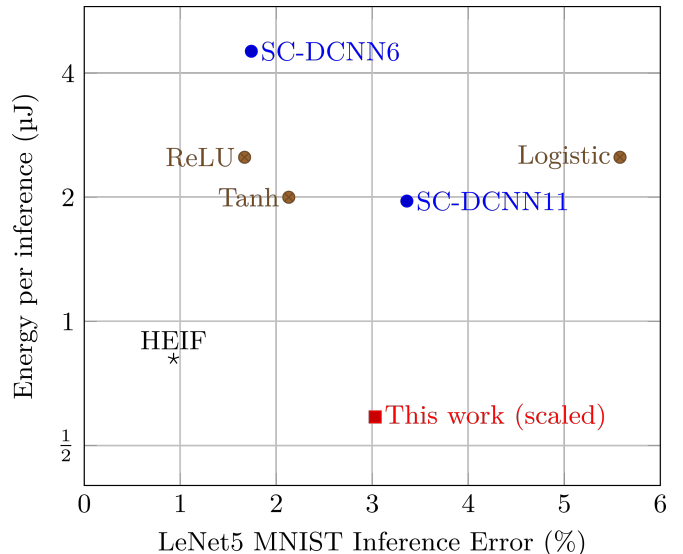


Figure 9. Comparison of stochastic computing implementations of LeNet5 on the MNIST dataset from the literature. The survey of logistic, ReLU, and tanh networks is reported in Ref. [69]. Two best-performing examples were extracted for SC-DCNN [65], and HEIF is reported in Ref. [66]. These references are all reported at the 45 nm technology node. The results for our work are presented at 1 V with $N = 128$ and $\delta = 16$ and are scaled up by a factor of four to bring our 22 nm node calculation into fair comparison with the literature.

therefore capacitance) to make a fair comparison [84]. We found that with significant energy savings we could achieve comparable accuracy, 97 %. While HEIF, the most modern work in that field, does significantly better (around 99.1 %), it is also the culmination of a research program that optimized all the other networks in Fig. 9, which originally performed at higher energy and lower performance. We have made no similar attempt to optimize our network and believe that it could be improved in principle to perform with similar accuracy and significantly lower energy.

The question remains as to whether 97 % is a useful recognition accuracy. The answer is conditionally yes, as long as one chooses the correct application. Though our network addressed the “Hello, world!” task of handwriting recognition, one might imagine a similar convolutional neural network used for face recognition. Presumably, our implementation would not perform as well as what can be done in mainstream neural networks. For mission critical applications like biometric identification, this would be an unacceptable drop in performance. But for, say, automatic face detection in a power-constrained edge context like a mobile device’s camera application, the drop in performance may be well worth the considerable energy efficiency advantages.

VI. CONCLUSION

We introduce a hybrid approach to classical stochastic computing, based on truly stochastic, low energy bitstreams generated by SMTJs. We introduce energy efficient primitive circuit elements (SC-PCsAs) that can be used to interface SMTJs with standard stochastic circuits. To test their effectiveness and explore the relaxed design space constraints afforded by true randomness, we simulated a neural network with OR gate neurons driven by SMTJ-based bitstream generators. Stochastic simulations of the network give 97 % classification accuracy on the MNIST dataset and circuit simulations of the circuit elements show that the energy usage should be about 150 nJ per inference, several times less than other stochastic implementations of LeNet5.

The energy efficiency we find in this case study is made possible in large part by true randomness, which eliminates the source correlations that would arise from the use of periodic pseudorandom number generators. Timing-based decorrelators called isolators address but do not entirely eliminate graph correlations, which arise due to reconvergent fanout of stochastic bitstreams. We analyze the energy and power usage of our neural network and discuss, in particular and in general, how SMTJ-based stochastic computing will scale with the progress being made in materials science. From the neural network perspective, our case study of LeNet5 was not optimized over all potential hyperparameters, so we expect that it should be possible to achieve higher classification accuracy and greater energy efficiency with more research. Non-neuromorphic stochastic computing designs [85] may also be able to benefit from correlation-free randomness sources. The design space degrees of freedom enabled by true randomness are an exciting landscape for future research.

ACKNOWLEDGMENTS

MWD and AMa contributed equally to this work. We thank Julie Grollier, Keven Garello, William Borders, Timothy Sherwood, George Tzimpragos, and Brian Hoskins for enlightening discussions and comments on the manuscript. MWD, AMa, PT, and AMi acknowledge support under the Cooperative Research Agreement Award No. 70NANB14H209, through the University of Maryland.

Appendix A: The SR-latch in the SC-PSCA

Judiciously adding the set-reset (SR) latch is a crucial modification we make to the standard PSCA circuit. The SR latch is a standard electrical engineering concept that can be found in textbooks [86]. Because the operating principles of the SR latch may nevertheless be unfamiliar to some readers, we include in this appendix

a brief description of the activity in the circuit in Fig. 1. Typical descriptions of the SR-latch are given in terms of two recurrently connected NAND gates, but we will analyze it at the transistor level since the operation of a field-effect transistor is more physically transparent and familiar to many physicists.

We described in Sec. II how nodes e and f are brought to either 01 or 10 after the clock signal goes to 1, the so-called evaluation phase. It seems that the state of node f is a good representative of the current state of the SMTJ, but unfortunately node f goes to a high voltage state in the pre-charge phase when the clock goes back to 0. Therefore f carries artifacts of the circuit operation, and does not faithfully represent the physical state of the SMTJ at the time it was previously read. We fix this by adding the SR-latch in Fig. 1(b).

First, consider the conclusion of the evaluation phase for the case $e = 1$ and $f = 0$, as illustrated in Fig. 10(a,b). In this figure, blue wires are at V_{dd} while red wires are connected to ground. The evaluated voltages at e and f are physically wired to the e and f branches of the latch in Fig. 10(b). Because branch e is on, the top-left transistor path is disconnected, denoted by a red circle. One can verify by hand that both transistors above the output node are disconnected (in the sense of open switches), while the transistors below the output node are connected to ground. The transistors above and below node g are in the opposite configuration. The output node thus reliably captures the state of node f during the evaluation phase. Notice that there are no paths for current to flow from V_{dd} to ground, so ohmic losses terminate as soon as this steady state is reached.

Now we consider the pre-charge phase in Fig. 10(c,d). Both nodes e and f are brought to a high voltage state in preparation for the RC race that will occur when the clock goes high again. Now that node f is high, the top and bottom transistors change state on the f branch of the latch, in the far right of Fig. 10(d). However, this does not change the state of the output node (or node g), node g (and the output node) themselves contribute to the opening and closing of each other's paths to V_{dd} and ground. Again, there is no steady state current flow once equilibrium is reached, so there are no continuous ohmic losses.

The two parallel transistors above node g act like a OR gate; if either of their gate voltages is low, then a current path exists between V_{dd} and node g . The series transistors below node g act like an AND gate; a current path exists from g to ground only if both gate voltages are high. This logical analysis will lead one to the standard presentation, but for our purposes it is sufficient to see that a change in state of f does not cause a change of state in the output node, so the unwanted comb structure from Fig. 2(d) is avoided.

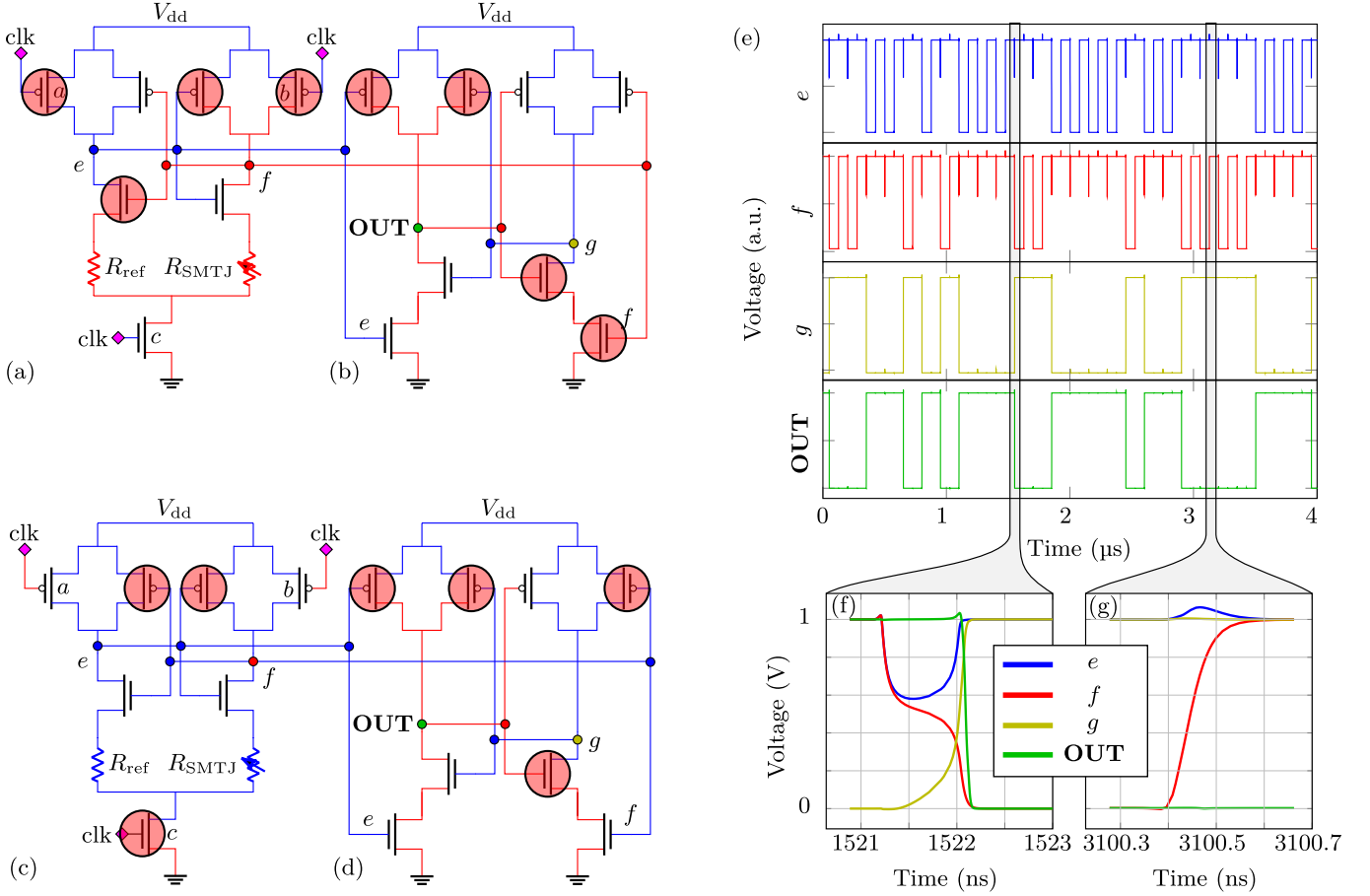


Figure 10. Operation of SC-PCSA. Blue wires represent $V = V_{dd}$, red wires represent $V = 0$, and red circles indicate that the source-drain path of transistor is turned off. (a): steady-state conclusion of evaluation phase, corresponding to the far-right of (f). The PCSA has detected that the SMTJ is in the parallel state, pulling node f to ground and node e to V_{dd} as a result. (b): The SR-latch copies the value of node f to the output node. (c): steady-state conclusion of reset/pre-charge, corresponding to the far-right of (g). (d): The PCSA has brought both e and f to V_{dd} in preparation for the measurement in the next evaluation phase. Although f has been brought to V_{dd} , the internal state of the latch does not change, so the output node continues to correctly represent the previously measured state of the SMTJ. (e, f, g): States of the various circuit node voltages throughout a simulation of the SC-PCSA. Curve colors correspond to points in the circuit indicated by the small circles in (a-d).

Appendix B: Comparison with stochastically-switched MTJs

A popular method in the literature for random number generation based on magnetic nanotechnology is the use of stochastic switching phenomena in *non-volatile* MTJs [31, 44, 52]. In this approach, an MTJ is first written to one preferred configuration, and then a sub-critical write voltage is placed across an MTJ. The strength of the write voltage is related to the parameter p of the Bernoulli trial to be performed. After probabilistically writing to the MTJ, the state is sensed with a read voltage in the usual way; the MTJ will have switched to the other stable configuration with probability p .

To compare our method with this one, we refer to Ref. [44] as a standard point of reference. Though that reference does not report numbers for the entire circuitry, they do report resistance and voltage numbers for the

tunnel junction itself. Their devices have an average resistance $R = 1000 \, \Omega$, and they apply an average of 217.5 mV bias perturbative write voltages over an average window of 8.75 ns. We neglect the cost of their read operation and any required CMOS circuitry here; the ohmic losses in the MTJ alone amount to 414 fJ per bit. Ref. [52] uses similar techniques on MTJs specifically tailored to the generation of stochastic computing bitstreams (the same task we consider here); they report average costs of 526 fJ per bit.

Though the energy per bit is orders of magnitude higher than what we have presented in Fig. 5, the tasks are not immediately equivalent. The main difference seems to be in speed available in the current state-of-the-art. Write/read cycles on non-volatile MTJs can be quite fast, allowing for rapid generation of random numbers in principle. Our scheme and that of [32] avoid large ohmic losses by using sense amplifiers, but as a result are limited

in speed by the natural time scale of the SMTJ's random telegraph noise. Experiments on SMTJs have demonstrated millisecond scale fluctuations; uncorrelated random bits cannot be sampled significantly faster than this time scale, which goes as roughly $e^{\Delta/k_B T}$ ns. If faster SMTJs can be engineered, our scheme would continue to function at the same energy per bit performance that we describe in this paper, and could plausibly produce bits down at the nanosecond time scale for ultra-low barrier devices.

Appendix C: Comparison with p -bits

A common application of SMTJs is the p -bit [14–17, 19–21]. A p -bit consists of an SMTJ and a transistor in series between a supply voltage V_{dd} and ground. The node between the SMTJ and the transistor is a voltage divider; its output is sent to a sequence of specially tuned CMOS inverters, where that output fluctuates around some reference voltage to produce a random telegraph noise signal. Controlling the gate-source voltage of the transistor in an analog fashion allows one to access analog tunability of the spin-torque on the SMTJ; this consequently allows one to tune the probability encoded in the output of the CMOS inverters.

Generally speaking, ideal p -bits are in a different class of device than what we propose here. In particular, they operate on analog voltage inputs. One could in principle build an analog probabilistic computer out of these units, whereas our proposed circuit is purely digital. Here, we do not compare the relative computational ability of digital versus analog computing but evaluate how well p -bit circuits would work for stochastic computing.

To make a fair comparison of device power and energy, we assume the supply voltage is again $V_{dd} = 1$ V. The voltage division provided by the inverter and SMTJ in the p -bit swings between two values which must be distinguished by the CMOS inverter. In order to center the $p = 1/2$ response of the p -bit at an input voltage of $V = 0.5$ V, the effective resistance of the transistor at that input voltage should be $\sqrt{R_P R_{AP}}$ to maximize the sensitivity of the circuit to both the high and low resistance states. In that configuration, we essentially have a 1 V voltage drop across a series resistor of effective resistance $\sqrt{R_P R_{AP}} + (R_P + R_{AP})/2$. The current running through the structure is therefore

$$I = \frac{2V_{dd}}{(R_P + R_{AP}) + 2\sqrt{R_P R_{AP}}}. \quad (C1)$$

For $R_{AP} = 100$ k Ω and $R_P = 50$ k Ω , this amounts to a 3.4 μ A current, or 3.4 μ W of ohmic dissipation per device. Neglecting the energy cost associated with generating independently controllable analog voltage sources for each p -bit, the ohmic loss is about 500 fJ per bit at the 150 ns clock cycle we consider, substantially greater than the ≈ 10 fJ for our PCSA-based approach. However,

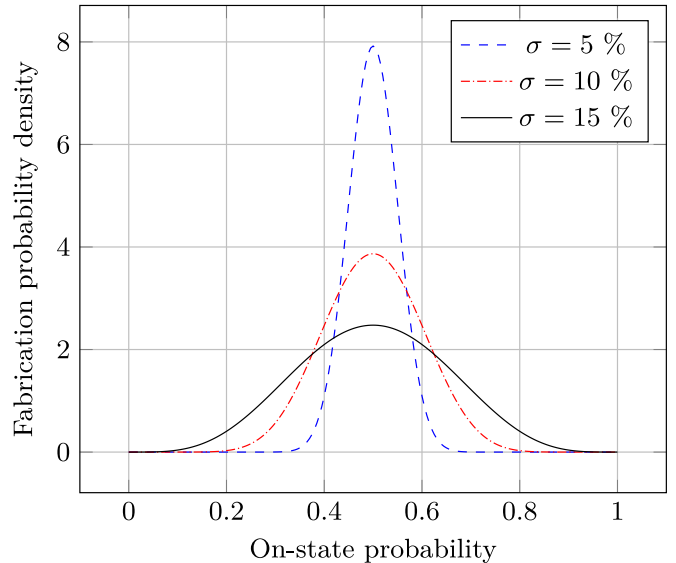


Figure 11. Probability density functions for the beta distribution centered at $1/2$ with three different standard deviations.

as the clock speed of the circuit is increased, the PCSA-based approach has constant energy cost per bit while the p -bit approach scales roughly linearly to lower energies with increasing speed. In order to achieve 10 fJ per bit performance, the autocorrelation time of the stochastic fluctuations in a p -bit would need to be reduced to 3 ns.

Recent work has attempted to push the theory of SMTJs toward this limit where the analog behavior can be harnessed more energy efficiently. Ref. [18] makes similar order of magnitude energy projections as ours in the case that the SMTJ dwell time could be reduced to ≈ 1 ns. Theory [40] suggests that nanomagnets with autocorrelation times on this scale might be realizable in the limit that the barrier goes to zero. Realizing devices that operate in this regime face a number of obstacles. As the barrier goes to zero, the current becomes more and more efficient at dictating the SMTJ state, requiring a delicate balance between adequate read currents and currents that control the state of the nanomagnet. Fabricating such a device may require extremely narrow margins of error.

Appendix D: Device variability

In the main text, especially leading up to Fig. 6, we assume a distribution for device variability. However, we do not assume the distribution to be gaussian, because the exponential tails of the normal distribution are problematic; the probability that a particular device is found in the parallel or antiparallel states is *strictly* confined to the open subset $(0, 1)$, by construction. Instead of trying to truncate or morph normal distributions to the domain, we opt instead to use the beta distribution, which is nat-

urally defined on the unit interval.

Suppose that the probability for a device to be in the on-state is itself a random variable \mathcal{P} . We choose to model \mathcal{P} as a beta distribution with unnormalized probability density

$$f(p) \propto \begin{cases} p^{\mu\phi-1}(1-p)^{\phi-\mu\phi-1} & 0 \leq p \leq 1 \\ 0 & \text{otherwise.} \end{cases} \quad (\text{D1})$$

The normalization is given by an Euler integral of the first kind. A convenient parametrization of the beta distribution is to specify its mean μ and a shape parameter $\phi > 0$, which is related to the variance of the distribution by

$$\sigma^2 = \frac{\mu(1-\mu)}{1+\phi}. \quad (\text{D2})$$

First, consider the case where $\mu = 1/2$ but ϕ is greater than zero. We plot three different distributions with standard deviations $\sigma \in \{0.05, 0.1, 0.15\}$ in Fig. 11. The existence of these distributions at the inputs to the bitstream generator will induce distributions in the statistics of the generator’s output bitstream. Expressions for these output distributions are analytically tractable [87], but are given in terms of hypergeometric Meijer G-functions and provide little useful intuition. To gain insight into the relationship between device variance and bitstream variance, we sample five thousand generator outputs at each programming code, for each of the distributions in Fig. 11. Fig. 6 plots these distributions vertically at each programmable probability.

A similar idea for producing a discretely-specified probability bitstream from a collection of LFSRs, called the weighted binary generator (WBG), has been discussed in Refs. [88, 89]. Our circuit has the advantage of $O(n)$ scaling of the input capacitance as a function of bit resolution n , whereas the input capacitance of the WBG scales as $O(n^2)$. We expect our solution to take less energy in general, regardless of whether LFSRs or SMTJs are used as the randomness sources. On the other hand, our numerical experiments indicate that for large deviations of the mean input probabilities from $1/2$, distortions in the output behavior of the WBG are more well-behaved than ours in the sense that the ordering of outputs is preserved relative to the WBG’s binary programming. In the case that programming protocols can be determined after fabricated devices are characterized, our circuit remains favorable. But if the mean of beta distribution governing expected probability for each SMTJ is unknown, or if the SMTJ variance is very large and programming protocol needs to be uniform across many devices, then the use of the WBG circuit may be preferable for reliability and ease of programming.

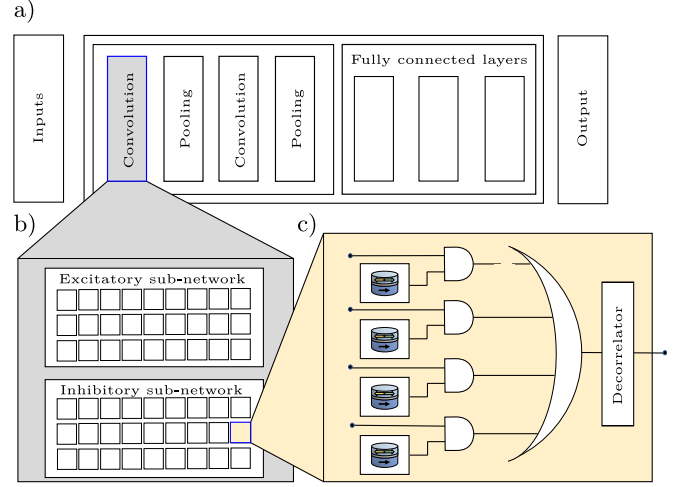


Figure 12. (a) Schematic block diagram of LeNet5 neural network. (b) Schematic block diagram of a neural layer (either a convolutional or a fully connected layer). Each of these layers are composed by an excitatory and inhibitory sub-networks. (c) Schematic view of a functional neural block comprised of superparamagnetic tunnel junctions, AND gates, one OR gate, and a decorrelator.

Appendix E: Details of the neural network

1. Layer structures

Each layer in a stochastic neural network architecture receives output bitstreams from the previous layer, as well as bitstreams from the programmable SMTJ weight arrays. These bitstreams are all multiplied in parallel using AND gates, the outputs of which are fed into the OR gate neurons for summation and activation. In other words, a fully connected layer with input degree m and output degree n is simply implemented as n different m -input OR gate neurons as we described in Sec. IV. The output of each layer is then passed through a random mask decorrelator as we described in Appendix E 4.

In a convolutional layer, we spatially multiplex the kernel applications. This means that each weight, which is being generated as described in Sec. II, is fanned out to a large number of input pixels at the same time. This massive weight sharing will cause neighboring pixels in the layer output to be strongly correlated, making the use of a graph decorrelator especially important here.

Today, most neural network designers implement the max operation as their pooling function, though average-pooling and min-pooling are also used. Ref. [90] has developed a correlation-insensitive max function for stochastic circuits that could be used for this purpose. Our pooling layer performs “OR-pooling”—that is, we simply use a single 4-input OR gate as our 2×2 pooling operation. This is similar to the method used in Ref. [91] to approximate average pooling, but without their weighting elements.

2. Dual architecture

In this paper, we use unipolar encoding to represent values in the stochastic circuit. To use unipolar encoding as such is to identify the probability of a wire being in the on state with the value that wire is said to encode. A disadvantage of this approach is that only numbers between zero and one can be encoded by the network. An alternative is bipolar encoding, where a wire turned on with probability p is said to encode the value $2p - 1$ [92]. This approach is used by many others in the field [65, 66], but unfortunately it is not clear to us that the OR gate can still be used as a useful nonlinear activation function in this encoding scheme.

Inhibitory behavior is generally believed to be a crucial aspect of neural networks [93, 94]. In most previous stochastic network proposals, use of the bipolar representation [92] grants the network access to negative numbers and, consequently, a mechanism for inhibitory signals. Some work has also been done on learning in networks with nonnegative weights [95], where inhibition is exercised by amplifying all signals except for those targeted for inhibition. Unfortunately, the unipolar representation employed here can represent neither negative numbers nor numbers greater than unity. Unipolar stochastic networks therefore lack the usual inhibitory mechanisms and fail to meaningfully learn most nontrivial data sets in our experiments.

To address this issue, we use a variation on a method proposed in Refs. [73, 76]. We refer to our strategy as a *dual architecture*. We employ two separate weight matrices in each layer, which are labeled as the excitatory and inhibitory subnetworks. The entire output vector from the inhibitory subnetwork is then elementwise inverted and AND'd with the excitatory output to achieve the response of the full layer. In the absence of correlations, the j^{th} element z_j of the final output from a network layer is therefore given by

$$z_j = y_j^{(e)} \left(1 - y_j^{(i)} \right), \quad (\text{E1})$$

where $y_j^{(e,i)}$ is the output element for the excitatory (inhibitory) subnetwork. This leads more explicitly to Eq. (2) from the main text. In reality, the preliminary outputs $y_j^{(e)}$ and $y_j^{(i)}$ are never uncorrelated, because they are both sourced from the same set of inputs $\{x_j\}$ to the layer. However, the use of the approximation in Eq. (E1) is necessary to use standard, local backpropagation methods, and we find that any induced error is minimal.

3. Simulation and Training

We build two kinds of models for our architecture: an analytic one and a stochastic one. The analytic model provides a representation of the network in terms of probabilities rather than the more complicated stochastic representation of random processes, allowing it to run faster

than the stochastic model. The analytic model verifies the functional correctness of our network structure and provides a baseline against which to compare our stochastic simulations. The stochastic model is a logic-level simulator, which allows us to determine the effects of correlation and to estimate activity factors needed for energy estimates. It uses input bitstreams that are generated from statistics consistent with the known properties of SMTJs.

To train our network, we use the standard backpropagation algorithm applied to our nonstandard network functionality given by Eq. (2) (that is, Eqs. (1) and (E1)). The speed of the analytic model makes it suitable for the training phase, as running detailed stochastic simulations for each inference of the training process would be prohibitive. This model is deterministic in the sense that we use the analytic probability equations for inference and backpropagation; it is local in the sense that we assume the inputs to each node of the network to be free of correlations. To train the model, we use a mini-batch version of the RMSProp optimization algorithm [96]. We use a learning rate of 0.005 and a forgetting factor of 0.95. We train sixty randomly initialized models for sixteen epochs each and select the best performing model. This network is then transferred to the stochastic model, to be simulated using realistic bitstreams generated according to known statistics for physical SMTJ devices.

4. Decorrelation

In order to train our neural network efficiently, we use backpropagation and gradient descent. In practice, backpropagation is implemented as a local learning rule; the gradient of the cost function is determined at each node in the network as a function only of that node and its nearest, connected neighbors. This locality is essential in keeping backpropagation algorithmically efficient, and therefore necessitates the use of a local analytic description that assumes all inputs are statistically independent. Such an analytic description, however, will necessarily fail to capture most graph correlations. In the specific case of a neural network, then, we need additional functionality for addressing graph correlations.

Because of the massive fan-in required in neural network systems, the algorithm from Ref. [48] would be unwieldy to implement in a neural network. For a fully connected layer with n neurons on the output, total decorrelation would require different delays for each neuron, leading to $O(n^2)$ delay elements per neuron. Our approach is to instead delay each neuron by a pseudorandom but fixed amount, creating a random mask of delay lengths on the output of each neural network layer. Such a programmable feature is implemented in our architecture by a chain of flip flops which are tapped into a multiplexer. The output of each decorrelator is fed to the next layer. In other words, every output of every layer is delayed by some integer between zero and a fixed up-

per bound, chosen pseudorandomly and uniformly over the interval at the network programming step. This approach does not formally eliminate graph correlations, but it empirically reduces their impact to a low enough

level such that implementation of the network architecture becomes feasible. The results presented in Sec. V demonstrate that the delay time can be set long enough to mostly saturate correlation-based network errors.

-
- [1] D. Apalkov, B. Dieny, and J. Slaughter, Magnetoresistive random access memory, *Proceedings of the IEEE* **104**, 1796 (2016).
 - [2] T. Hanyu, T. Endoh, D. Suzuki, H. Koike, Y. Ma, N. Onizawa, M. Natsui, S. Ikeda, and H. Ohno, Standby-power-free integrated circuits using mtj-based vlsi computing, *Proceedings of the IEEE* **104**, 1844 (2016).
 - [3] Q. Guo, X. Guo, R. Patel, E. Ipek, and E. G. Friedman, Ac-dimm: associative computing with stt-mram, *ACM SIGARCH Computer Architecture News* **41**, 189 (2013).
 - [4] M. Bapna, S. K. Piotrowski, S. D. Oberdick, M. Li, C.-L. Chien, and S. A. Majetich, Magnetostatic effects on switching in small magnetic tunnel junctions, *Applied Physics Letters* **108**, 022406 (2016).
 - [5] A. Mizrahi, T. Hirtzlin, A. Fukushima, H. Kubota, S. Yuasa, J. Grollier, and D. Querlioz, Neural-like computing with populations of superparamagnetic basis functions, *Nature Communications* **9**, 1533 (2018).
 - [6] W. Rippard, R. Heindl, M. Pufall, S. Russek, and A. Kos, Thermal relaxation rates of magnetic nanoparticles in the presence of magnetic fields and spin-transfer effects, *Physical Review B* **84**, 064439 (2011).
 - [7] J. Sun and D. Ralph, Magnetoresistance and spin-transfer torque in magnetic tunnel junctions, *Journal of Magnetism and Magnetic Materials* **320**, 1227 (2008).
 - [8] S.-W. Lee and K.-J. Lee, Emerging three-terminal magnetic memory devices, *Proceedings of the IEEE* **104**, 1831 (2016).
 - [9] H. Sato, S. Ikeda, and H. Ohno, Magnetic tunnel junctions with perpendicular easy axis at junction diameter of less than 20 nm, *Japanese Journal of Applied Physics* **56**, 0802A6 (2017).
 - [10] N. Perrissin, G. Grégoire, S. Lequeux, L. Tillie, N. Strelkov, S. Auffret, L. Buda-Prejbeanu, R. Sousa, L. Vila, B. Dieny, *et al.*, Perpendicular shape anisotropy spin transfer torque magnetic random-access memory: towards sub-10 nm devices, *Journal of Physics D: Applied Physics* **52**, 234001 (2019).
 - [11] J. Grollier, D. Querlioz, and M. D. Stiles, Spintronic nanodevices for bioinspired computing, *Proceedings of the IEEE* **104**, 2024 (2016).
 - [12] K. Roy, A. Sengupta, and Y. Shim, Perspective: Stochastic magnetic devices for cognitive computing, *Journal of Applied Physics* **123**, 210901 (2018).
 - [13] K. Y. Camsari, B. M. Sutton, and S. Datta, p -bits for probabilistic spin logic, *Applied Physics Reviews* **6**, 011305 (2019).
 - [14] R. Faria, K. Y. Camsari, and S. Datta, Implementing bayesian networks with embedded stochastic mram, *AIP Advances* **8**, 045101 (2018).
 - [15] X. Jia, J. Yang, Z. Wang, Y. Chen, H. H. Li, and W. Zhao, Spintronics based stochastic computing for efficient bayesian inference system, in *2018 23rd Asia and South Pacific Design Automation Conference (ASP-DAC)* (IEEE, 2018) pp. 580–585.
 - [16] A. Sengupta, M. Parsa, B. Han, and K. Roy, Probabilistic deep spiking neural systems enabled by magnetic tunnel junction, *IEEE Transactions on Electron Devices* **63**, 2963 (2016).
 - [17] K. Y. Camsari, R. Faria, B. M. Sutton, and S. Datta, Stochastic p -bits for invertible logic, *Physical Review X* **7**, 031014 (2017).
 - [18] W. A. Borders, A. Z. Pervaiz, S. Fukami, K. Y. Camsari, H. Ohno, and S. Datta, Integer factorization using stochastic magnetic tunnel junctions, *Nature* **573**, 390 (2019).
 - [19] S. Ganguly, K. Y. Camsari, and A. W. Ghosh, Reservoir computing using stochastic p -bits, *arXiv preprint arXiv:1709.10211* (2017).
 - [20] O. Hassan, K. Y. Camsari, and S. Datta, Voltage-driven building block for hardware belief networks, *IEEE Design & Test* **36**, 15 (2019).
 - [21] Y. Shim, A. Jaiswal, and K. Roy, Ising computation based combinatorial optimization using spin-hall effect (she) induced stochastic magnetization reversal, *Journal of Applied Physics* **121**, 193902 (2017).
 - [22] A. Sengupta, P. Panda, P. Wijesinghe, Y. Kim, and K. Roy, Magnetic tunnel junction mimics stochastic cortical spiking neurons, *Scientific Reports* **6**, 30039 (2016).
 - [23] A. F. Vincent, J. Larroque, N. Locatelli, N. B. Romdhane, O. Bichler, C. Gamrat, W. S. Zhao, J.-O. Klein, S. Galdin-Retailleau, and D. Querlioz, Spin-transfer torque magnetic memory as a stochastic memristive synapse for neuromorphic systems, *IEEE transactions on biomedical circuits and systems* **9**, 166 (2015).
 - [24] G. Srinivasan, A. Sengupta, and K. Roy, Magnetic tunnel junction based long-term short-term stochastic synapse for a spiking neural network with on-chip stdp learning, *Scientific Reports* **6**, 29545 (2016).
 - [25] W. J. Gross and V. C. Gaudet, *Stochastic Computing: Techniques and Applications* (Springer, 2019).
 - [26] A. Alaghi, W. Qian, and J. P. Hayes, The promise and challenge of stochastic computing, *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* **37**, 1515 (2018).
 - [27] W. Qian, X. Li, M. D. Riedel, K. Bazargan, and D. J. Lilja, An architecture for fault-tolerant computation with stochastic logic, *IEEE Transactions on Computers* **60**, 93 (2010).
 - [28] V. T. Lee, A. Alaghi, R. Pamula, V. S. Sathe, L. Ceze, and M. Oskin, Architecture considerations for stochastic computing accelerators, *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* **37**, 2277 (2018).
 - [29] F. Neugebauer, I. Polian, and J. P. Hayes, S-box-based random number generation for stochastic computing, *Microprocessors and Microsystems* **61**, 316 (2018).
 - [30] J. M. de Aguiar and S. P. Khatri, Exploring the viability of stochastic computing, in *2015 33rd IEEE International Conference on Computer Design (ICCD)* (IEEE, 2015)

- pp. 391–394.
- [31] A. Fukushima, T. Seki, K. Yakushiji, H. Kubota, H. Imamura, S. Yuasa, and K. Ando, Spin dice: A scalable truly random number generator based on spintronics, *Applied Physics Express* **7**, 083001 (2014).
 - [32] D. Vodenicarevic, N. Locatelli, A. Mizrahi, J. S. Friedman, A. F. Vincent, M. Romera, A. Fukushima, K. Yakushiji, H. Kubota, S. Yuasa, S. Tiwari, J. Grollier, and D. Querlioz, Low-energy truly random number generation with superparamagnetic tunnel junctions for unconventional computing, *Physical Review Applied* **8**, 054045 (2017).
 - [33] B. Parks, M. Bapna, J. Igbokwe, H. Almasi, W. Wang, and S. A. Majetich, Superparamagnetic perpendicular magnetic tunnel junctions for true random number generators, *AIP Advances* **8**, 055903 (2018).
 - [34] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner, *et al.*, Gradient-based learning applied to document recognition, *Proceedings of the IEEE* **86**, 2278 (1998).
 - [35] W. Zhao, C. Chappert, V. Javerliac, and J.-P. Noziere, High speed, high stability and low power sensing amplifier for mtj/cmos hybrid logic circuits, *IEEE Transactions on Magnetics* **45**, 3784 (2009).
 - [36] W. Zhao and Y. Cao, New generation of predictive technology model for sub-45 nm early design exploration, *IEEE Transactions on Electron Devices* **53**, 2816 (2006).
 - [37] A. PTM, Arizona state university predictive technology model, 22nm, <http://ptm.asu.edu/>.
 - [38] A. Mizrahi, N. Locatelli, R. Matsumoto, A. Fukushima, H. Kubota, S. Yuasa, V. Cros, J. Kim, J. Grollier, and D. Querlioz, Magnetic stochastic oscillators: Noise-induced synchronization to underthreshold excitation and comprehensive compact model, *IEEE Transactions on Magnetics* **51**, 1 (2015).
 - [39] R. Faria, K. Y. Camsari, and S. Datta, Low-barrier nanomagnets as *p*-bits for spin logic, *IEEE Magnetics Letters* **8**, 1 (2017).
 - [40] J. Kaiser, A. Rustagi, K. Y. Camsari, J. Z. Sun, S. Datta, and P. Upadhyaya, Ultrafast fluctuations in low-barrier magnets, *arXiv preprint arXiv:1902.03312* (2019).
 - [41] K. Tsuchida, T. Inaba, K. Fujita, Y. Ueda, T. Shimizu, Y. Asao, T. Kajiyama, M. Iwayama, K. Sugiura, S. Ikegawa, *et al.*, A 64mb mram with clamped-reference and adequate-reference schemes, in *2010 IEEE International Solid-State Circuits Conference (ISSCC)* (IEEE, 2010) pp. 258–259.
 - [42] Y. Jiang, T. Nozaki, S. Abe, T. Ochiai, A. Hirohata, N. Tezuka, and K. Inomata, Substantial reduction of critical current for magnetization switching in an exchange-biased spin valve, *Nature Materials* **3**, 361 (2004).
 - [43] O. Hassan, R. Faria, K. Y. Camsari, J. Z. Sun, and S. Datta, Low-barrier magnet design for efficient hardware binary stochastic neurons, *IEEE Magnetics Letters* **10**, 1 (2019).
 - [44] W. H. Choi, Y. Lv, J. Kim, A. Deshpande, G. Kang, J.-P. Wang, and C. H. Kim, A magnetic tunnel junction based true random number generator with conditional perturb and real-time output probability tracking, in *2014 IEEE International Electron Devices Meeting* (IEEE, 2014) pp. 12–5.
 - [45] A. Rukhin, J. Soto, J. Nechvatal, M. Smid, and E. Barker, *A statistical test suite for random and pseudorandom number generators for cryptographic applications*, Tech. Rep. (Booz-Allen and Hamilton Inc Mclean Va, 2001).
 - [46] D. Jenson and M. Riedel, A deterministic approach to stochastic computation, in *2016 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)* (IEEE, 2016) pp. 1–8.
 - [47] M. H. Najafi, D. J. Lilja, and M. Riedel, Deterministic methods for stochastic computing using low-discrepancy sequences, in *2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)* (IEEE, 2018) pp. 1–8.
 - [48] P.-S. Ting and J. P. Hayes, Isolation-based decorrelation of stochastic circuits, *2016 IEEE 34th International Conference on Computer Design (ICCD)* (2016).
 - [49] G. E. Salam and R. M. Goodman, A digital neural network architecture using random pulse trains, in *Silicon Implementation of Pulse Coded Neural Networks*, edited by M. E. Zaghloul, J. L. Meador, and R. W. Newcomb (Springer US, Boston, MA, 1994) pp. 249–261.
 - [50] H. Ichihara, S. Ishii, D. Sunamori, T. Iwagaki, and T. Inoue, Compact and accurate stochastic circuits with shared random number sources, in *2014 IEEE 32nd International Conference on Computer Design (ICCD)* (IEEE, 2014) pp. 361–366.
 - [51] K. Kim, J. Lee, and K. Choi, An energy-efficient random number generator for stochastic circuits, in *2016 21st Asia and South Pacific Design Automation Conference (ASP-DAC)* (2016).
 - [52] A. Mondal and A. Srivastava, Energy-efficient design of mtj-based neural networks with stochastic computing, *J. Emerg. Technol. Comput. Syst.* **16**, 7:1 (2019).
 - [53] S.-W. Chung, T. Kishi, J. Park, M. Yoshikawa, K. Park, T. Nagase, K. Sunouchi, H. Kanaya, G. Kim, K. Noma, *et al.*, 4gbit density stt-mram using perpendicular mtj realized with compact cell structure, in *2016 IEEE International Electron Devices Meeting (IEDM)* (IEEE, 2016) pp. 27–1.
 - [54] K. Rho, K. Tsuchida, D. Kim, Y. Shirai, J. Bae, T. Inaba, H. Noro, H. Moon, S. Chung, K. Sunouchi, *et al.*, 23.5 a 4gb lpddr2 stt-mram with compact 9f2 1t1mtj cell and hierarchical bitline architecture, in *2017 IEEE International Solid-State Circuits Conference (ISSCC)* (IEEE, 2017) pp. 396–397.
 - [55] N. Perrissin, S. Lequeux, N. Strelkov, A. Chavent, L. Vila, L. D. Buda-Prejbeanu, S. Auffret, R. C. Sousa, I. L. Prejbeanu, and B. Dieny, A highly thermally stable sub-20 nm magnetic random-access memory based on perpendicular shape anisotropy, *Nanoscale* **10**, 12187 (2018).
 - [56] E. Liu, J. Swerts, Y. C. Wu, A. Vaysset, S. Couet, S. Mertens, S. Rao, W. Kim, S. Van Elshocht, J. De Boeck, *et al.*, Top-pinned stt-mram devices with high thermal stability hybrid free layers for high-density memory applications, *IEEE Transactions on Magnetics* **54**, 1 (2018).
 - [57] K. Tsunoda, M. Aoki, H. Noshiro, Y. Iba, S. Fukuda, C. Yoshida, Y. Yamazaki, A. Takahashi, A. Hatada, C. Nakabayashi, *et al.*, Area dependence of thermal stability factor in perpendicular stt-mram analyzed by bi-directional data flipping model, in *2014 IEEE International Electron Devices Meeting* (IEEE, 2014) pp. 19–3.
 - [58] C. Park, H. Lee, C. Ching, J. Ahn, R. Wang, M. Pakala, and S. Kang, Low ra magnetic tunnel junction arrays in conjunction with low switching current and high breakdown voltage for stt-mram at 10 nm and beyond, in *2018*

- IEEE Symposium on VLSI Technology* (IEEE, 2018) pp. 185–186.
- [59] L. Xue, C. Ching, A. Kontos, J. Ahn, X. Wang, R. Whig, H.-w. Tseng, J. Howarth, S. Hassan, H. Chen, *et al.*, Process optimization of perpendicular magnetic tunnel junction arrays for last-level cache beyond 7 nm node, in *2018 IEEE Symposium on VLSI Technology* (IEEE, 2018) pp. 117–118.
 - [60] J. G. Alzate, U. Arslan, P. Bai, J. Brockman, Y. J. Chen, N. Das, K. Fischer, T. Ghani, P. Heil, P. Hentges, R. Jahan, A. Littlejohn, M. Mainuddin, D. Ouellette, J. Pellegrin, T. Pramanik, C. Puls, P. Quintero, T. Rahman, M. Sekhar, B. Sell, M. Seth, A. J. Smith, A. K. Smith, L. Wei, C. Wiegand, O. Golonzka, and F. Hamzaoglu, 2.4 2 mb array-level demonstration of stt-mram process and performance towards l4 cache applications, in *2019 IEEE International Electron Devices Meeting (IEDM)* (IEEE, 2019).
 - [61] W. S. McCulloch and W. Pitts, A logical calculus of the ideas immanent in nervous activity, *The bulletin of mathematical biophysics* **5**, 115 (1943).
 - [62] J. Von Neumann, Probabilistic logics and the synthesis of reliable organisms from unreliable components, (1956).
 - [63] W. Poppelbaum, C. Afuso, and J. Esch, Stochastic computing elements and systems, in *Proceedings of the November 14-16, 1967, Fall Joint Computer Conference* (ACM, 1967) pp. 635–644.
 - [64] B. R. Gaines, Origins of stochastic computing, in *Stochastic Computing: Techniques and Applications* (Springer, 2019) pp. 13–37.
 - [65] A. Ren, Z. Li, C. Ding, Q. Qiu, Y. Wang, J. Li, X. Qian, and B. Yuan, Sc-dcnn: highly-scalable deep convolutional neural network using stochastic computing, *Proceedings of the Twenty-Second International Conference on Architectural Support for Programming Languages and Operating Systems – ASPLOS '17* (2017).
 - [66] Z. Li, J. Li, A. Ren, R. Cai, C. Ding, X. Qian, J. Draper, B. Yuan, J. Tang, and Q. Qiu, Heif: highly efficient stochastic computing based inference framework for deep neural networks, *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* (2018).
 - [67] A. Ardakani, F. Leduc-Primeau, N. Onizawa, T. Hanyu, and W. J. Gross, Vlsi implementation of deep neural network using integral stochastic computing, *IEEE Transactions on Very Large Scale Integration (VLSI) Systems* **25**, 2688 (2017).
 - [68] N. Onizawa, W. J. Gross, and T. Hanyu, Brain-inspired computing, in *Stochastic Computing: Techniques and Applications*, edited by W. J. Gross and V. C. Gaudet (Springer International Publishing, 2019) pp. 185–199.
 - [69] J. Li, Z. Yuan, Z. Li, C. Ding, A. Ren, Q. Qiu, J. Draper, and Y. Wang, Hardware-driven nonlinear activation for stochastic computing based deep convolutional neural networks, in *2017 International Joint Conference on Neural Networks (IJCNN)* (IEEE, 2017) pp. 1230–1236.
 - [70] K. Kim, J. Lee, and K. Choi, An energy-efficient random number generator for stochastic circuits, in *2016 21st Asia and South Pacific Design Automation Conference (ASP-DAC)* (IEEE, 2016) pp. 256–261.
 - [71] K. Kim, J. Kim, J. Yu, J. Seo, J. Lee, and K. Choi, Dynamic energy-accuracy trade-off using stochastic computing in deep neural networks, in *Proceedings of the 53rd Annual Design Automation Conference, DAC '16* (ACM, New York, NY, USA, 2016) pp. 124:1–124:6.
 - [72] A. Alaghi and J. P. Hayes, Exploiting correlation in stochastic circuit design, in *2013 IEEE 31st International Conference on Computer Design (ICCD)* (IEEE, 2013) pp. 39–46.
 - [73] J. Tomlinson, M.S., D. Walker, and M. Sivilotti, A digital neural network architecture for vlsi, 1990 IJCNN International Joint Conference on Neural Networks (1990).
 - [74] Y.-C. Kim and M. Shanblatt, An implementable digital multilayer neural network (dmnn), 1992 IJCNN International Joint Conference on Neural Networks (1992).
 - [75] Y.-C. Kim and M. Shanblatt, Random noise effects in pulse-mode digital multilayer neural networks, *IEEE Transactions on Neural Networks* **6**, 220 (1995).
 - [76] Y.-C. Kim and M. Shanblatt, Architecture and statistical model of a pulse-mode digital multilayer neural network, *IEEE Transactions on Neural Networks* **6**, 1109 (1995).
 - [77] V. Sze, Y.-H. Chen, T.-J. Yang, and J. S. Emer, Efficient processing of deep neural networks: A tutorial and survey, *Proceedings of the IEEE* **105**, 2295 (2017).
 - [78] S. Gupta, A. Agrawal, K. Gopalakrishnan, and P. Narayanan, Deep learning with limited numerical precision, in *International Conference on Machine Learning* (2015) pp. 1737–1746.
 - [79] S. Zhou, Y. Wu, Z. Ni, X. Zhou, H. Wen, and Y. Zou, Dorefa-net: training low bitwidth convolutional neural networks with low bitwidth gradients, *arXiv preprint arXiv:1606.06160* (2016).
 - [80] I. Hubara, M. Courbariaux, D. Soudry, R. El-Yaniv, and Y. Bengio, Quantized neural networks: Training neural networks with low precision weights and activations, *The Journal of Machine Learning Research* **18**, 6869 (2017).
 - [81] C. Song, B. Liu, W. Wen, H. Li, and Y. Chen, A quantization-aware regularized learning method in multi-level memristor-based neuromorphic computing system, in *2017 IEEE 6th Non-Volatile Memory Systems and Applications Symposium (NVMSA)* (IEEE, 2017) pp. 1–6.
 - [82] Q. Yang, H. Li, and Q. Wu, A quantized training method to enhance accuracy of rram-based neuromorphic systems, in *2018 IEEE International Symposium on Circuits and Systems (ISCAS)* (IEEE, 2018) pp. 1–5.
 - [83] R. Manohar, Comparing stochastic and deterministic computing, *IEEE Computer Architecture Letters* **14**, 119 (2015).
 - [84] J.-A. Carballo, W.-T. J. Chan, P. A. Gargini, A. B. Kahng, and S. Nath, Itrs 2.0: Toward a re-framing of the semiconductor technology roadmap, in *2014 IEEE 32nd International Conference on Computer Design (ICCD)* (IEEE, 2014) pp. 139–146.
 - [85] A. Alaghi and J. P. Hayes, Strauss: Spectral transform use in stochastic circuit synthesis, *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems* **34**, 1770 (2015).
 - [86] N. H. E. Weste and K. Eshraghian, *Principles of VLSI Design: A Systems Perspective*, 2nd ed. (Addison-Wesley, 1994).
 - [87] M. D. Springer and W. Thompson, The distribution of products of beta, gamma and gaussian random variables, *SIAM Journal on Applied Mathematics* **18**, 721 (1970).
 - [88] P. K. Gupta and R. Kumaresan, Binary multiplication with pn sequences, *IEEE Transactions on Acoustics, Speech, and Signal Processing* **36**, 603 (1988).
 - [89] M. Yang, B. Li, D. J. Lilja, B. Yuan, and W. Qian, Towards theoretical cost limit of stochastic number generators for stochastic computing, in *2018 IEEE Computer*

- Society Annual Symposium on VLSI (ISVLSI)* (IEEE, 2018) pp. 154–159.
- [90] F. Neugebauer, I. Polian, and J. P. Hayes, On the maximum function in stochastic computing, in *Proceedings of the 16th ACM International Conference on Computing Frontiers – CF ’19* (2019).
 - [91] B. Li, Y. Qin, B. Yuan, and D. J. Lilja, Neural network classifiers using stochastic computing with a hardware-oriented approximate activation function, in *2017 IEEE International Conference on Computer Design (ICCD)* (IEEE, 2017) pp. 97–104.
 - [92] C. Winstead, Tutorial on stochastic computing, in *Stochastic Computing: Techniques and Applications* (Springer, 2019) pp. 39–76.
 - [93] D. R. Chialvo and P. Bak, Learning from mistakes, *Neuroscience* **90**, 1137 (1999).
 - [94] P. Bak and D. R. Chialvo, Adaptive learning by extremal dynamics and negative feedback, *Physical Review E* **63**, 031912 (2001).
 - [95] J. Chorowski and J. M. Zurada, Learning understandable neural networks with nonnegative weight constraints, *IEEE Transactions on Neural Networks and Learning Systems* **26**, 62 (2015).
 - [96] G. Hinton, N. Srivastava, and K. Swersky, Neural networks for machine learning, lecture 6a: Overview of mini-batch gradient descent (2012).