

Reliability and entropy production in nonequilibrium electronic memories

Peer-reviewed author version

Freitas, Nahuel; PROESMANS, Karel & Esposito, Massimiliano (2022) Reliability and entropy production in nonequilibrium electronic memories. In: PHYSICAL REVIEW E, 105 (3) (Art N° 034107).

DOI: 10.1103/PhysRevE.105.034107

Handle: <http://hdl.handle.net/1942/37164>

Reliability and entropy production in non-equilibrium electronic memories

Nahuel Freitas,¹ Karel Proesmans,^{1,2} and Massimiliano Esposito¹

¹*Complex Systems and Statistical Mechanics, Department of Physics and Materials Science, University of Luxembourg, L-1511 Luxembourg, Luxembourg*

²*Hasselt University, B-3590 Diepenbeek, Belgium*

(Dated: April 4, 2022)

We find the relation between reliability and entropy production in a realistic model of electronic memory (low-power MOS-based SRAM) where logical values are encoded as metastable non-equilibrium states. We employ large deviations techniques to obtain an analytical expression for the bistable quasipotential describing the non-equilibrium steady state and use it to derive an explicit expression bounding the error rate of the memory. Our results go beyond the dominant contribution given by classical instanton theory and provide accurate estimates of the error rate as confirmed by comparison with stochastic simulations.

PACS numbers:

I. INTRODUCTION

A common strategy to reduce the energy consumption of electronic computing devices is to reduce the voltage at which they are powered. However, this strategy is limited by the fact that as the operation voltage is reduced, different sources of electrical noise start to play an increasingly important role [1–4]. The most fundamental and unavoidable one is given by the thermal fluctuations intrinsic to any device. It originates from the interaction with degrees of freedom that are not explicitly described, but that can be normally assumed to be at thermal equilibrium. A rigorous description of intrinsic thermal noise in complex and non-linear electronic circuits is thus a fundamental problem in modern engineering, of great importance for the search of new efficient computing schemes [3–7]. However, it is also a hard problem that is usually given approximate treatments involving different kinds of approximations that are difficult to control, and that in general compromise thermodynamic consistency [8, 9]. This issue was recently addressed by the development of a general theoretical framework to construct thermodynamically consistent stochastic models of non-linear electronic circuits [9].

In this work we make use of that framework to analyze the tradeoff between reliability and dissipation (i.e. entropy production) of low-power static random access memory (SRAM) cells. Due to their speed and low energy consumption, SRAM cells are employed as internal memory in virtually all modern processors. The occurrence of errors induced by thermal noise in low-power implementations has been mainly studied using numerical methods based on stochastic simulations [4, 10]. The reason is that in low-power regimes current fluctuations are Poissonian and cannot be faithfully described as Gaussian noise [11], which considerably complicates analytical treatments. However, since one is typically interested in determining the rate of errors in regimes where errors are rare, the amount of computational time demanded by the stochastic simulations can be extremely large [4]. In this contribution we report two main results. First, we obtain

an analytical description of the steady state fluctuations of the memory, fully capturing the non-equilibrium transition from a monostable phase into the bistable phase that allows the representation of a bit. Secondly, we show how to employ the previous result to analytically estimate the error rate of the memory. By comparing with exact stochastic simulations, we show that our analytical estimation correctly describes the scaling of the error rate with the voltage that powers the memory. Then, we show that the error rate is exponentially suppressed as the square of the dissipation (for large dissipation). To get there, we make use of advanced methods from stochastic thermodynamics [9], large deviations theory [12–16], and first-passage time statistics [17–19].

II. BASIC MODEL

We consider the usual model of a SRAM memory cell core: two inverters, or NOT gates, connected in a loop (see Figure 1-(a)). In particular, we consider the implementation based on complementary metal-oxide-semiconductor (MOS) transistors. In this case, each inverter is itself composed of an n MOS transistor and a p MOS transistor. The circuit is powered by applying a voltage bias $\Delta V = V_{\text{dd}} - V_{\text{ss}}$. The deterministic and linear stability analysis of the circuit (see Appendix A) shows that for low values of ΔV the circuit has a unique fixed point, but when ΔV is above a critical value there is a transition into bistability, which is employed to encode a single bit of information. The transistors are modelled as externally controlled conduction channels with associated capacitances (see Figure 1-(b)). The charge conduction through each transistor channel is modelled as a bidirectional Poisson process. Thus, to each (n/p)MOS transistor we associate two Poisson rates $\lambda_{\pm}^{n/p}(V_{\text{GS}}, V_{\text{DS}})$, where the subindices \pm correspond to the forward and backward conduction directions, and V_{GS} and V_{DS} are the gate-source and drain-source voltage drops, respectively. For fixed voltages V_{dd} and V_{ss} the circuit has two independent degrees of freedom: the voltages v_1 and

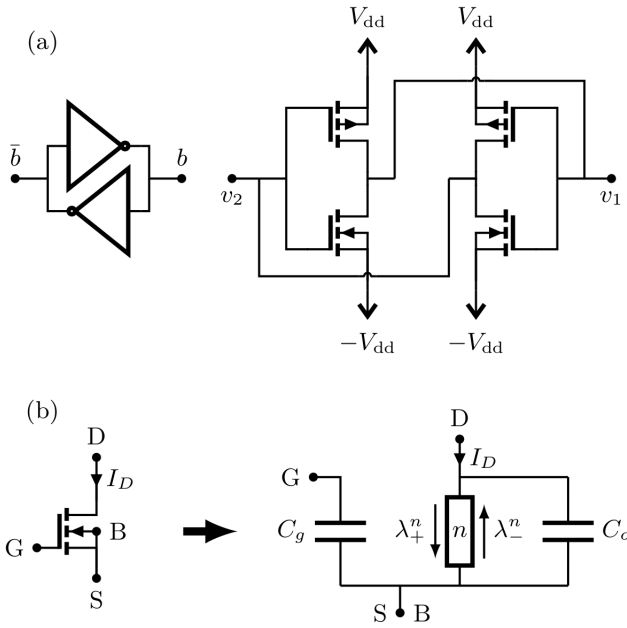


FIG. 1: (a) A bistable logical circuit constructed with two NOT gates, representing a bit, and its CMOS implementation, where each NOT gate is constructed with one p MOS (top) and one n MOS (bottom) transistors. (b) Each transistor (in this example an n MOS one) is modelled as a conduction channel between drain (D) and source (S) terminals, with associated rates λ_{\pm}^n . The gate-body (G-B) interface is represented as a capacitor C_g , and another capacitor C_o takes into account the output capacitance. Other parasitic capacitances could also be taken into account, for example between drain and gate. With this model and taking $V_{ss} = -V_{dd}$, the total electrostatic energy of the full circuit is $\Phi(v_1, v_2) = (C/2)(v_1^2 + v_2^2) + CV_{dd}^2$, with $C = 2(C_o + C_g)$.

v_2 at the outputs of each inverter. These are discrete stochastic quantities, that in principle can only take the values mv_e , where m is any integer and $v_e = q_e/C$ (q_e is the positive electron charge and C a value of capacitance characterizing the device, see Figure 1-(b)).

At any given time the state of the system is described by a probability distribution $P(v_1, v_2, t)$ over the state space. Its evolution is given by the following master equation

$$d_t P(v_1, v_2, t) = PA|_{v_1-v_e, v_2} + PB|_{v_1+v_e, v_2} + PA^*|_{v_1, v_2-v_e} + PB^*|_{v_1, v_2+v_e} - P(A + B + A^* + B^*)|_{v_1, v_2}, \quad (1)$$

where we are using the compact notation $PA|_{v_1, v_2} = P(v_1, v_2, t)A(v_1, v_2)$, and $A^*(v_1, v_2) = A(v_2, v_1)$. The transition rates $A(v_1, v_2)$ and $B(v_1, v_2)$ are combinations of the Poisson rates assigned to the transistors:

$$\begin{aligned} A(v_1, v_2) &= \lambda_+^p(v_1, v_2) + \lambda_-^n(v_1, v_2) \\ B(v_1, v_2) &= \lambda_-^p(v_1, v_2) + \lambda_+^n(v_1, v_2). \end{aligned} \quad (2)$$

In order to guarantee thermodynamic consistency, the

Poisson rates $\lambda_{\pm}^{n/p}(v_1, v_2)$ must satisfy the so called *local detailed balance (LDB)* conditions. As an example, for the p MOS transistor in the first inverter, this condition reads:

$$\frac{\lambda_+^p(v_1, v_2)}{\lambda_-^p(v_1 + v_e, v_2)} = e^{-\delta Q/(k_b T)}, \quad (3)$$

where $\delta Q = \Phi(v_1 + v_e, v_2) - \Phi(v_1, v_2) - q_e V_{dd}$, $\Phi(v_1, v_2)$ is the electrostatic energy of the system, and we have considered the environment of the transistor to be in equilibrium at temperature T . For $V_{ss} = -V_{dd}$, as we will consider in the following, the electrostatic energy is $\Phi(v_1, v_2) = (C/2)(v_1^2 + v_2^2) + CV_{dd}^2$. Thus the LDB condition of Eq. (1) relates the rates of the transitions $v_1 \rightleftharpoons v_1 + v_e$ to the difference in internal energy between those states, and the work $q_e V_{dd}$ realized by the voltage sources during the forward transition. Then, δQ is the total energy change associated to that transition, and since by energy conservation it must be provided by the environment of the device, it is the heat interchanged with it. A condition analogous to Eq. (3) is imposed to all the transistors present in the circuit. A general procedure to construct thermodynamically consistent rates based on the I-V curve characterization of a given devices was recently identified in [9]. For the case of MOS transistors in subthreshold operation, one obtains:

$$\begin{aligned} \lambda_+^p(v_1, v_2) &= (I_0/q_e) e^{(V_{dd}-v_2-V_{th})/(nV_T)} \\ \lambda_-^p(v_1, v_2) &= \lambda_+^p(v_1, v_2) e^{-(V_{dd}-v_1)/V_T} e^{-(v_e/2)/V_T}, \end{aligned} \quad (4)$$

and $\lambda_{\pm}^n(v_1, v_2) = \lambda_{\pm}^p(-v_1, -v_2)$. In the previous equation $V_T = k_b T/q_e$ is the thermal voltage and I_0 , V_{th} , and n are parameters characterizing the transistor (respectively known as *specific current*, *threshold voltage*, and *slope factor*). An incorrect procedure to construct transition rates, which is however used in some numerical simulations [4, 10], is to employ the rates directly obtained from the I-V curve characterization, without enforcing the LDB conditions. In that way one finds rates that are obtained from the ones of Eq. (4) by removing the factor $e^{-(v_e/2)/V_T}$ appearing in $\lambda_{\pm}^{p/n}$. Although this factor is in many situations very close to 1, it can become relevant for small devices or at low temperatures, and it is in fact responsible for the charging effects in single-electron devices [20–22]. Also, neglecting that factor leads to systematic errors in the determination of the steady state. For example, in modern CMOS fabrication processes capacitance values as low as $C \simeq 50$ aF can be attained [23], which correspond to elementary voltages as high as $v_e \simeq 3$ mV. At room temperature we have $V_T \simeq 26$ mV and therefore $v_e/V_T \simeq 0.1$ and $e^{-(v_e/2)/V_T} \simeq 0.95$.

For mathematical simplicity, the parameters I_0 , V_{th} and n are considered to be the same for all the four transistors involved in the circuit. That is, we are not taking into account any variability associated with the fabrication process [24]. It should be possible to extend our results to systems with asymmetric parameters.

III. STEADY STATE DISTRIBUTION AND LARGE DEVIATIONS PRINCIPLE

To find the steady state of the memory one option is to construct the generator of the master equation in Eq. (1) and compute its eigenvector of zero eigenvalue (see Figure 1-(a)). Analytical progress is possible by considering a macroscopic limit and employing the principle of large deviations. This limit consists in assuming that the elementary voltage v_e is negligible compared to all other voltage scales, which in this case are the thermal voltage V_T and the powering voltage V_{dd} (thus, the limit $v_e \rightarrow 0$ used in the following must be interpreted as $v_e/V_T \rightarrow 0$ and $v_e/V_{dd} \rightarrow 0$ for fixed V_T and V_{dd}). Physically, this corresponds to large devices, for which the typical capacitance C is large and thus v_e is small. Also, from Eq. (4) we have that the Poisson rates are proportional to $(I_0/C)v_e^{-1}$. As explained in Appendix B, the specific current I_0 can also be considered to be proportional to the size of the device, and therefore we see that the transition rates scale as v_e^{-1} . Under these conditions, as $v_e \rightarrow 0$, the deterministic equations of motion are recovered from the master equation in Eq. (1) (see Appendixes A and B), and one also expects the distribution $P_{ss}(v_1, v_2)$ to become strongly peaked around the deterministic stationary values [19, 25]. In this context, the LD principle states that departures from the deterministic values are suppressed exponentially in v_e^{-1} . This is expressed mathematically as the existence of the limit $f(v_1, v_2) = \lim_{v_e \rightarrow 0} -v_e \log(P_{ss}(v_1, v_2))$, or equivalently [14]:

$$P_{ss}(v_1, v_2) \asymp e^{-(f(v_1, v_2) + o(v_e))/v_e}. \quad (5)$$

Therefore, as $v_e \rightarrow 0$, the values of v_1 and v_2 will be perfectly localized at a global minimum of the *rate function* $f(v_1, v_2)$. Indeed, the minima of $f(v_1, v_2)$ correspond to the deterministic fixed points (see Appendixes A and B). We will refer to the function $f(v_1, v_2)$ as a *quasipotential* describing the steady state distribution. This is in analogy to an equilibrium situation, where the steady state must be the equilibrium Boltzmann distribution $P_{eq}(v_1, v_2) \propto \exp(-\Phi(v_1, v_2)/k_b T)$ and thus, by Eq. (5), $f(v_1, v_2)$ should match the true potential energy $\Phi(v_1, v_2)$ scaled by the thermal voltage V_T . Also, the interpretation of $f(v_1, v_2)$ as a potential has a deeper justification on the fact that it always is a Lyapunov function for the deterministic dynamics [25], as the true potential energy is for equilibrium settings.

Plugging Eq. (5) into Eq. (1), imposing $d_t P_{ss} = 0$, and only keeping the lower order terms in v_e , we obtain the following differential equation for $f(v_1, v_2)$:

$$0 = (e^{\partial_{v_1} f} - 1) a(v_1, v_2) + (e^{-\partial_{v_1} f} - 1) b(v_1, v_2) \\ + (e^{\partial_{v_2} f} - 1) a(v_2, v_1) + (e^{-\partial_{v_2} f} - 1) b(v_2, v_1), \quad (6)$$

where $a(v_1, v_2) = \lim_{v_e \rightarrow 0} v_e A(v_1, v_2)$, and the same for $b(v_1, v_2)$. The same equation can be obtained by more

general path integral methods, in terms of a Hamiltonian defining an action in the space of all possible stochastic trajectories [15, 25–27]. This equation cannot be solved exactly. However, an approximate solution can be found by exploiting the fact that the variables $x = (v_1 - v_2)/2$ and $y = (v_1 + v_2)/2$ are, except for some trivial correlations discussed below, approximately independent. Thus, as explained in Appendix B, from Eq. (6) the rate functions $g(x)$ and $h(y)$ corresponding to the partial distributions $P(x) = \sum_y P_{ss}(y + x, y - x) \asymp \exp(-g(x)/v_e)$ and $Q(y) = \sum_x P_{ss}(y + x, y - x) \asymp \exp(-h(y)/v_e)$ can be found to be

$$d_x g(x) = 2 \log \left(\frac{a(-x, 0) + b(x, 0)}{a(x, 0) + b(-x, 0)} \right), \\ d_y h(y) = 2 \log \left(\frac{b(x_{\min}, y) + b(-x_{\min}, y)}{a(x_{\min}, y) + a(-x_{\min}, y)} \right), \quad (7)$$

where the change of variables in the functions a and b is understood, and x_{\min} in the expression for $d_y h$ is the minimum of $g(x)$. The given expression for $d_x g(x)$ is actually exact, since it only relies on the fact that the most probable value of y for any x is always $y = 0$ in the limit $v_e \rightarrow 0$ (as can be seen from the symmetry of the exact steady state, see Figure 1-(a)), and does not require x and y to be considered independent variables.

The variables x and y will be always correlated because, since v_1/v_e and v_2/v_e are integer random variables, their difference $2x/v_e$ and sum $2y/v_e$ will always have the same parity. If, however, when restricted to a given parity, x and y can be considered independent, and if both parities have the same probability, then the full probability distribution $P_{ss}(v_1, v_2)$ can be reconstructed from the partial distributions $P(x)$ and $Q(y)$ as

$$P_{ss}(y + x, y - x) = 2P(x)Q(y)\text{Par}(x, y), \quad (8)$$

where $\text{Par}(x, y)$ is one if $2x/v_e$ and $2y/v_e$ have the same parity, or zero if they do not. Eq. (8) allows to approximately reconstruct the full steady state distribution from the partial rate functions $g(x)$ and $h(y)$. As shown in Appendix C, this approximation becomes exact for typical fluctuations in the low-noise regime $v_e/V_T \ll 1$, and is extremely accurate in general.

The results in Eq. (7) are in principle valid for any Poisson rates $\lambda_{\pm}^{n/p}$. Remarkably, for the particular MOS rates of Eq. (4), the expression for $d_x g$ can be integrated exactly, resulting in:

$$g(x) = \frac{x^2 + 2V_{dd}x}{V_T} + \frac{2nV_T}{n+2} [L(x, V_{dd}) - L(x, -V_{dd})], \quad (9)$$

where $L(x, V_{dd}) = \text{Li}_2(-\exp((V_{dd} + x(1 + 2/n))/V_T))$, and $\text{Li}_2(\cdot)$ is the polylogarithm function of second order. This is the first important result of this work, and will allow us to analytically estimate the error rate of a low-power SRAM memory cell in the next section. In turn, the rate function $h(y)$ can be seen to satisfy $h(y) = h_0 y^2/V_T + \mathcal{O}(y^4)$ (an expression for h_0 in terms of the circuit parameters is given in Appendix B).

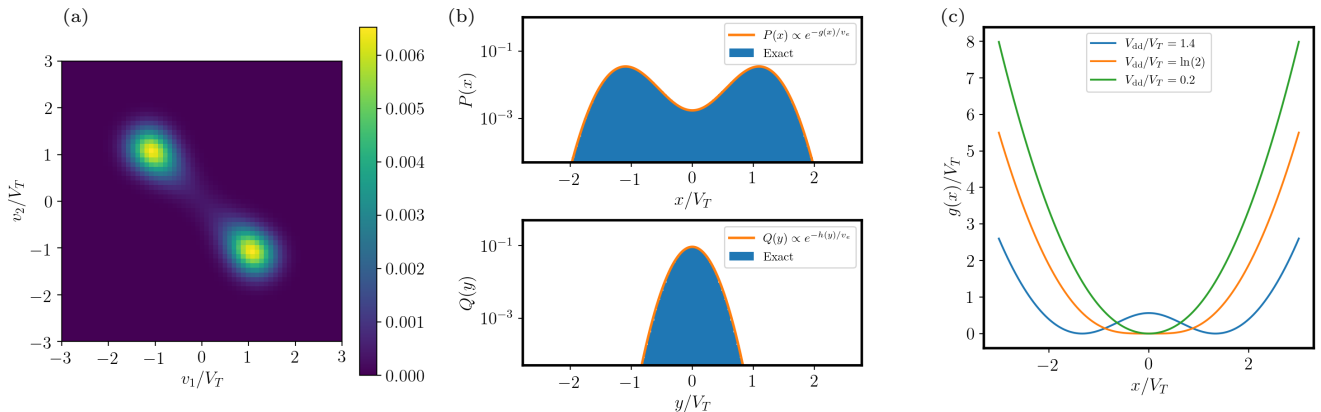


FIG. 2: (a) Exact steady state obtained by numerical integration of the master equation ($V_{\text{dd}}/V_T = 1.2$, $v_e/V_T = 0.1$, $n=1$). (b) Partial distributions for the variables x and y as obtained from the exact global distribution in (a), and from the analytical results of Eq. (7). (c) Quasipotential $g(x)/v_e$ for different values of V_{dd} ($v_e/V_T = 0.1$, $n=1$).

In Figure 2-(a) we show the exact steady state distribution $P_{\text{ss}}(v_1, v_2)$ obtained by numerically evolving Eq. (1) for $v_e/V_T = 0.1$, $V_{\text{dd}}/V_T = 1.2$, and $n = 1$. We see that for these parameters the most probable values are distributed around $v_1 = -v_2 \simeq \pm V_{\text{dd}}$, i.e., the possible solutions to the deterministic equations of motion (Appendix A). In Figure 2-(b) we compare the exact partial distributions $P(x)$ and $Q(y)$ for the variables $x = (v_1 - v_2)/2$ and $y = (v_1 + v_2)/2$, respectively, with the ones obtained from the quasipotentials $g(x)/v_e$ and $h(y)/v_e$. We see that the agreement is remarkable despite the value of v_e being only one order of magnitude lower than V_T and V_{dd} (states with only a few tens of electrons are occupied). Finally, in Figure 2-(c) we show the quasipotential $g(x)$ for different values of the powering voltage V_{dd} . We see that there is a transition between a unimodal steady state and the bimodal distribution compatible with bistability, that for $n = 1$ happens at $V_{\text{dd}} = \ln(2)V_T$ (the data-retention voltage), as can also be seen from the analysis of the deterministic equations (Appendix A).

IV. ERROR RATE

If the initial state of the system is close to one of the two possible metastable NESSs, let us say $v_1 = -v_2 \simeq V_{\text{dd}}$, the ensuing dynamics will be characterized by two different time scales. First, a fast relaxation on the local basin of attraction will take place. Indeed, from the deterministic equations (Appendix A) we see that this relaxation develops at a rate $\lambda_{\text{eq}} \simeq \tau_0^{-1}(v_e/V_T) e^{2V_{\text{dd}}/V_T}$ that increases exponentially with V_{dd} , where $\tau_0 = (q_e/I_0) e^{V_{\text{th}}/(nV_T)}$ is a natural time scale for this problem. After this local metastable NESS has been reached, a slow dynamics consisting of rare transitions to the other possible metastable NESS follows. Since the metastable NESSs are associated to the values

of the stored bit, this rare transitions are considered errors. We are interested in computing the error rate λ_{err} in terms of the circuit parameters. This is a hard problem that has been mainly treated numerically [4, 10], and for which a rigorous stochastic treatment is crucial. It is possible to see that, to leading order in v_e^{-1} , the rate of escape out of a NESSs centered around $\mathbf{v}_{\text{min}} = (v_1^{\text{min}}, v_2^{\text{min}})$ can be obtained from the quasipotential f thanks to the following result [25, 28]:

$$\lim_{v_e \rightarrow 0} v_e \log(\tau_0 \lambda_{\text{err}}) = -(f(\mathbf{v}^*) - f(\mathbf{v}_{\text{min}})), \quad (10)$$

where \mathbf{v}^* is a saddle point of the quasipotential (which in this case is $\mathbf{v}^* = (0, 0)$). The factor $\exp(-(f(\mathbf{v}^*) - f(\mathbf{v}_{\text{min}}))/v_e)$ is also the dominant contribution to the probability of a trajectory, or ‘instanton’, going from \mathbf{v}_{min} to \mathbf{v}^* [25]. This result can be considered a generalization to NESSs of the classical Arrhenius’s law [29], and in this case leads to the ‘dominant’ estimate of the error rate

$$\lambda_{\text{err}}^{\text{D}} = \tau_0^{-1} e^{-(g(0) - g(x_{\text{min}}))/v_e}, \quad (11)$$

that can be readily evaluated from Eq. (9). However, this estimate misses any contribution to λ_{err} that is subexponential in v_e^{-1} , but that might be anyway relevant for finite values of v_e . For equilibrium systems some subexponential factors are provided by the classic Eyring-Kramers formula [30–32], in terms of the curvature of the energy surface at the fixed and saddle points. For out of equilibrium systems with Gaussian noise, subexponential corrections are discussed in [28, 33]. In our case, since we are dealing with a discrete out of equilibrium system subjected to shot noise, we resort to the general method explained in the following.

The first step to compute λ_{err} is to provide an operational definition of what an error is. We consider that the state of the memory is read by monitoring the output of

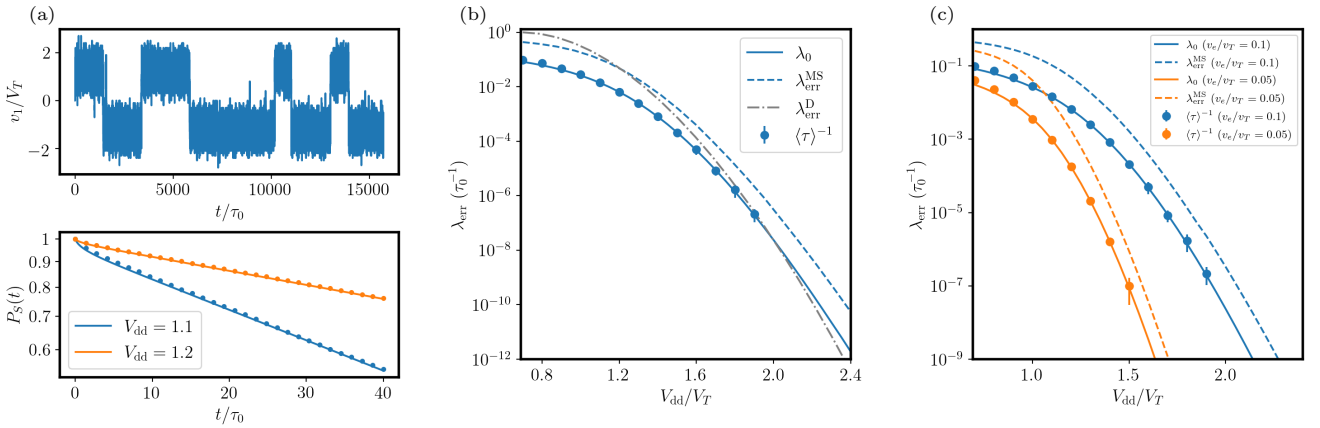


FIG. 3: (a) Sample trajectory generated with the Gillespie simulation of the stochastic dynamics (top, $V_{dd}/V_T = 1.2$, $v_e/V_T = 0.1$, $n=1$), and decay of the survival probability $P_S(t)$ for the protocol described in the text, for different values of V_{dd} (bottom). Solid lines were obtained by Eq. (12) and the dots from data generated with the Gillespie algorithm ($v_e/V_T = 0.1$, $n=1$). (b) Different estimates of the error rate as a function of V_{dd} for $v_e/V_T = 0.1$ and $n=1$. The dots indicate the inverse of the mean TTE, $\langle \tau \rangle^{-1}$, as obtained from Gillespie simulations. The solid blue line corresponds to the minimum eigenvalue λ_0 of the partial generator $-W_{\text{HH}}$, and the violet line to the metastable rate $\lambda_{\text{err}}^{\text{MS}}$ of Eq. (14). The dashed grey line shows the dominant contribution in the $v_e/V_T \rightarrow 0$ limit of Eq. (11). (c) Estimates of the error rate as a function of V_{dd} for different values of v_e/V_T ($n=1$).

the first inverter, i.e., the voltage v_1 . A zero or positive value of v_1 is identified with the logical state H ('high'), and a negative value with the logical state L ('low'). This logical encoding induces natural projection operations in the state space, that we construct as follows. Each microscopic state (v_1, v_2) is mapped to a vector $|v_1, v_2\rangle$. A given probability distribution $P(v_1, v_2)$ is represented as the vector $|P\rangle = \sum_{v_1, v_2} P(v_1, v_2) |v_1, v_2\rangle$, while the generator of the master equation in Eq. (1) is represented as a matrix \mathbb{W} acting over these vectors. Thus, the steady state distribution $|P_{\text{ss}}\rangle$ satisfies $0 = \mathbb{W}|P_{\text{ss}}\rangle$. The orthogonal projectors corresponding to the logical states H and L are, respectively, $\Pi_H = \sum_{v_1 \geq 0, v_2} |v_1, v_2\rangle \langle v_1, v_2|$ and $\Pi_L = \sum_{v_1 < 0, v_2} |v_1, v_2\rangle \langle v_1, v_2|$ (where $\langle a|$ is just the transpose of $|a\rangle$). Note that $\Pi_j \Pi_k = \delta_{j,k} \Pi_j$ and that $\Pi_H + \Pi_L = \mathbb{1}$. Then, we can consider the projections of the steady state to each of the logical subspaces: $|P_{\text{ss}}^H\rangle = \Pi_H |P_{\text{ss}}\rangle / \langle 1 | \Pi_H | P_{\text{ss}} \rangle$ and $|P_{\text{ss}}^L\rangle = \Pi_L |P_{\text{ss}}\rangle / \langle 1 | \Pi_L | P_{\text{ss}} \rangle$ ($|1\rangle$ is just the vector with unit components). Now we give the following operational definition of an error: at time $t = 0$ we prepare the system at a state drawn from the metastable distribution $|P_{\text{ss}}^H\rangle$ (for which the voltage v_1 is always positive or zero), and monitor its evolution until v_1 becomes negative. This event is considered an error, and the random time τ at which it takes place is recorded. We are interested in the distribution of τ , which can be considered a *first-passage* problem [18, 19]. As explained in [19], one possible approach to obtain the statistics of τ is to consider an alternative dynamics with absorbing boundary conditions at the interface between the logical subspaces. Thus, the *survival probability* of not observing

any error up to time t is given by

$$P_S(t) = \langle 1 | e^{\mathbb{W}_{\text{HH}} t} | P_{\text{ss}}^H \rangle. \quad (12)$$

Here, the matrix \mathbb{W}_{HH} is the partial generator $\Pi_H \mathbb{W} \Pi_H$ reduced to the H -subspace. The vectors $|1\rangle$ and $|P_{\text{ss}}^H\rangle$ are also reduced to the same subspace. The probability to observe an error between times t and $t + dt$ is $p(t)dt$, where $p(t) = -d_t P_S(t)$. Then the average time to an error (TTE) is

$$\langle \tau \rangle = \int_0^\infty \tau p(\tau) d\tau = \int_0^\infty P_S(\tau) d\tau. \quad (13)$$

At variance with the full generator \mathbb{W} , the partial generator \mathbb{W}_{HH} does not conserve probability (since it continuously leaks into the L -subspace), and therefore its largest eigenvalue is strictly lower than 0. Indeed, we can write $P_S(t) = \sum_k C_k e^{-\lambda_k t}$, where $-\lambda_k$ are the eigenvalues of \mathbb{W}_{HH} (with $0 < \lambda_0 \leq \lambda_1 \leq \lambda_2 \leq \dots$), and C_k are constants that depend on the initial state (with $\sum_k C_k = 1$). Thus, for large times we have $P_S(t) \simeq C_0 \exp(-\lambda_0 t)$. From this, it follows that for long times the distribution of τ is approximately exponential with rate λ_0 . This already provides a method to estimate the error rate: one should construct the generator \mathbb{W}_{HH} and numerically compute the eigenvalue of smallest absolute value, which can be done efficiently with several routines since the matrix \mathbb{W}_{HH} is sparse. Note that λ_0 is independent of the initial state. It is possible to obtain analytically another estimate of the error rate by exploiting the metastability of the initial state $|P_{\text{ss}}^H\rangle$. For this, we consider an approximation in which the state $|P(t)\rangle = e^{\mathbb{W}_{\text{HH}} t} |P_{\text{ss}}^H\rangle$ evolving according to the generator \mathbb{W}_{HH} is assumed to be always

proportional to $|P_{ss}^H\rangle$ (the initial distribution), but with a time dependent normalization. In that case the survival probability satisfies $d_t P_S(t) = \langle 1 | \mathbb{W}_{HH} | P_{ss}^H \rangle P_S(t)$ and therefore we can write $P_S(t) = e^{-\lambda_{err}^{MS} t}$, with the ‘metastable’ rate $\lambda_{err}^{MS} = -\langle 1 | \mathbb{W}_{HH} | P_{ss}^H \rangle$. This is equivalent to assume that the error rate is constant and equal to the initial one, and consequently depends explicitly on the initial state. Note that by the conservation of probability of the full generator ($\langle 1 | \mathbb{W} = 0$), and the property $\Pi_L + \Pi_H = \mathbb{1}$, we have the alternative expression $\lambda_{err}^{MS} = \langle 1 | \mathbb{W}_{LH} | P_{ss}^H \rangle$, where \mathbb{W}_{LH} is the reduction of the matrix $\Pi_L \mathbb{W} \Pi_H$ to the appropriate subspaces. This last expression for λ_{err}^{MS} can be evaluated using Eq. (8) for the steady state, with the following result:

$$\lambda_{err}^{MS} = 4 \sum_{v_2} B(0, v_2) P(-v_2/2) Q(v_2/2), \quad (14)$$

where $B(v_1, v_2)$ is given in Eq. (2), and $P(x)$ and $Q(y)$ are the LD approximations to the partial distributions, i.e., $P(x) \propto \exp(-g(x)/v_e)$ and $Q(y) \propto \exp(-h(y)/v_e)$. It is instructive to see how Eq. (14) reduces to Eq. (11) for $v_e \rightarrow 0$. First, we notice that $Q(y)$ becomes strongly peaked around $y = 0$ for $v_e \rightarrow 0$, and therefore we can approximate $\lambda_{err}^{MS} \simeq 4B(0, 0)P(0)$. In turn, we have $P(0) = \exp(-g(0)/v_e)/N$ with $N = \sum_x \exp(-g(x)/v_e)$, that for $v_e \rightarrow 0$ becomes $N \simeq \exp(-g(x_{min})/v_e)$. Then, we recover the result of Eq. (11), with τ_0^{-1} replaced by the factor $4B(0, 0)$ (which is subexponential, since the rates scale as v_e^{-1}).

In general there is no definite relation between the estimates λ_0 and λ_{err}^{MS} , and the mean TTE $\langle \tau \rangle$. However, for the particular protocol we are considering, in which the initial state is $|P_{ss}^H\rangle$, the instantaneous decay rate of the survival probability $\lambda(t) = -d_t \log(P_S(t))$ is a monotonously decreasing function. This is easily understood: the steady state distribution has a non-zero value at the boundary $v_1 = 0$ between logical subspaces. Then, the initial occupation of the states at or close to the boundary will quickly leak into the L -subspace, with a rate that decreases as the occupation of those states decrease, reaching its asymptotic value λ_0 for long times. In that case, from Eq. (13) it follows that the inverse of the average TTE is bounded by λ_0 and λ_{err}^{MS} :

$$\lambda_0 \leq \langle \tau \rangle^{-1} \leq \lambda_{err}^{MS}. \quad (15)$$

Thus, λ_{err}^{MS} provides an upper bound to the inverse mean TTE.

In Figure 3-(a) we show a sample trajectory obtained by the Gillespie algorithm, and the decay of the survival probability, computed with two methods. The solid lines were obtained from Eq. (12), by constructing the reduced generator \mathbb{W}_{HH} . The dots were obtained from Gillespie simulations in which initial states were drawn from the steady state distribution and the time to an error was recorded. We see that the decay rate decreases monotonously from the initial one to the asymptotic one given by λ_0 . From the same data we compute the mean

TTE $\langle \tau \rangle$. In Figure 3-(b) we compare $\langle \tau \rangle^{-1}$ with the different estimates of the error rate, as a function of V_{dd} . We see that λ_0 is an excellent estimate of $\langle \tau \rangle^{-1}$. The metastable rate λ_{err}^{MS} of Eq. (14) consistently overestimate the true error rate, but displays the same scaling with V_{dd} . In contrast, we see that the dominant estimate of Eq. (11) largely overestimate the error rate for low V_{dd} , while it underestimate it for large values of V_{dd} . Figure 3-(c) shows $\langle \tau \rangle^{-1}$, λ_0 and λ_{err}^{MS} as a function of V_{dd} for different values of v_e .

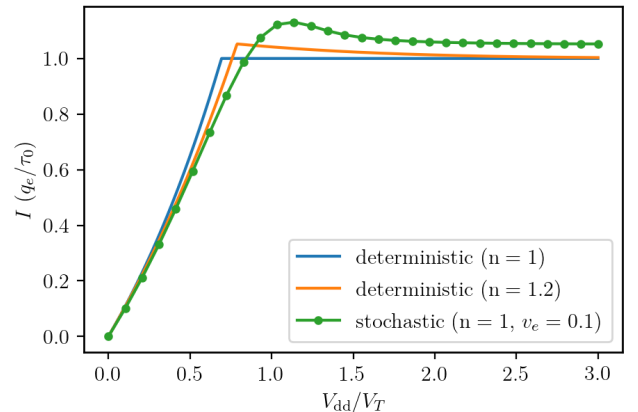


FIG. 4: Electrical current through each transistor in steady state conditions as a function of V_{dd} .

V. ENTROPY PRODUCTION

We now study the steady state entropy production of the memory. At steady state, the average current I through both inverters is the same. Thus, the rate at which heat is dissipated in the environment is $\dot{Q} = 4V_{dd}I$, and the entropy production rate is just $\dot{\Sigma} = \dot{Q}/T$. From the deterministic solution for $n = 1$, it follows that in the monostable phase the electric current increases exponentially with V_{dd} , $I = (q_e/\tau_0)(e^{V_{dd}/V_T} - 1)$, while it is constant in the bistable phase, $I = q_e/\tau_0$ (see Appendix A). In Figure 4 we show that the same constant value is achieved also for $n \neq 1$. In addition, we show the average current obtained by computing the mean value of $I(v_1, v_2) = q_e(\lambda_+^p(v_1, v_2) - \lambda_-^p(v_1, v_2))$ using the exact steady state distribution $P_{ss}(v_1, v_2)$. This average current also reaches a constant value for large V_{dd} , that is above the deterministic one due to finite- v_e effects. Interestingly, it displays a bump right after the onset of bistability. The origin of this maximum in the average current is precisely the occurrence of errors, since each switching event in which the memory flips its state has an associated dissipation. As V_{dd} increases, errors become rare and the average current tends to the value corresponding to any of the metastable NESSs with a

definite logical value.

Thus, for large V_{dd} the electrical current I is just constant, and therefore the entropy production $\dot{\Sigma}$ is proportional to V_{dd} . Also, from Eq. (9) it is possible to see that, to dominant order in $V_{\text{dd}} \gg V_T$, $\Delta g = g(0) - g(x_{\text{min}}) \simeq (2/(n+2))(V_{\text{dd}}^2/V_T)$. Then, it follows that for large entropy production rates the error rate scales as:

$$\begin{aligned} \lambda_{\text{err}}^{\text{MS}} &\propto e^{-\frac{2}{n+2} \frac{V_{\text{dd}}^2}{v_e V_T}} \\ &= e^{-\frac{2}{n+2} \frac{k_b T}{(4I)^2/C} (\dot{\Sigma}/k_b)^2}. \end{aligned} \quad (16)$$

Here we have ignored terms in $\log(\lambda_{\text{err}})$ that are constant or linear in $\dot{\Sigma}$ or equivalently V_{dd} , that can be easily included. Indeed, the previous equation is compatible with what was obtained in ad-hoc treatments based on Gaussian noise [1], up to model-dependent constant factors in the exponent. However, in general one must employ the result in Eq. (14), that can be readily evaluated.

VI. DISCUSSION

We used the theory of stochastic thermodynamics to construct a thermodynamically consistent stochastic model of a technologically relevant kind of electronic memory, subjected to Poissonian thermal noise. Large

deviations theory was then employed to obtain an analytical expression for the steady state of the memory, that allowed to estimate the rate at which errors occur. We have thus explicitly solved a problem that has been so far only treated using expensive numerical simulations [4].

From a wider perspective, our work shows how modern developments in statistical physics can contribute to solve important problems in electronic engineering. Although our focus has been on the problem of memory reliability, our methods and results are also relevant for the design of non-conventional stochastic computing schemes, where naturally occurring thermal fluctuations are exploited as a resource [5, 9, 34–36]. For instance, we note that our results directly apply to the low-power binary stochastic neuron proposed in [9] which is based on a SRAM memory cell core identical to the one studied here.

VII. ACKNOWLEDGMENTS

We acknowledge funding from the European Research Council, project NanoThermo (ERC-2015-CoGAgreement No. 681456), and from the Luxembourg National Research Fund (FNR), CORE project NTEC (C19/MS/13664907).

-
- [1] K. Natori and N. Sano, *Journal of applied physics* **83**, 5019 (1998).
 - [2] A. Wang, B. H. Calhoun, and A. P. Chandrakasan, *Sub-threshold design for ultra low-power systems*, vol. 95 (Springer, 2006).
 - [3] S. Krishnan, S. V. Garimella, G. M. Chrysler, and R. V. Mahajan, *IEEE Transactions on advanced packaging* **30**, 462 (2007).
 - [4] E. Rezaei, M. Donato, W. R. Patterson, A. Zaslavsky, and R. I. Bahar, *IEEE Transactions on Device and Materials Reliability* **20**, 488 (2020).
 - [5] J. Han and M. Orshansky, in *2013 18th IEEE European Test Symposium (ETS)* (IEEE, 2013), pp. 1–6.
 - [6] J. Gu and P. Gaspard, *Physical Review E* **99**, 012137 (2019).
 - [7] J. Gu and P. Gaspard, *Journal of Statistical Mechanics: Theory and Experiment* **2020**, 103206 (2020).
 - [8] P. Hänggi and P. Jung, *IBM Journal of Research and Development* **32**, 119 (1988).
 - [9] N. Freitas, J.-C. Delvenne, and M. Esposito, *Stochastic Thermodynamics of Non-Linear Electronic Circuits: A Realistic Framework for Thermodynamics of Computation* (2020), 2008.10578, URL <https://arxiv.org/abs/2008.10578v3>.
 - [10] H. Li, J. Mundy, W. Patterson, D. Kazazis, A. Zaslavsky, and R. Bahar, in *Proceedings of Workshop on System Effects of Logic Soft Errors* (2006).
 - [11] R. Sarpeshkar, T. Delbruck, and C. A. Mead, *IEEE Circuits and Devices Magazine* **9**, 23 (1993).
 - [12] V. Elgart and A. Kamenev, *Physical Review E* **70**, 041106 (2004).
 - [13] A. Kamenev and B. Meerson, *Physical Review E* **77**, 061107 (2008).
 - [14] H. Touchette, *Physics Reports* **478**, 1 (2009).
 - [15] M. Assaf and B. Meerson, *Journal of Physics A: Mathematical and Theoretical* **50**, 263001 (2017).
 - [16] D. T. Limmer, C. Y. Gao, and A. R. Poggioli, *arXiv* (2021), 2104.05194, URL <https://arxiv.org/abs/2104.05194v2>.
 - [17] P. Hänggi, P. Talkner, and M. Borkovec, *Reviews of modern physics* **62**, 251 (1990).
 - [18] S. Redner, *A guide to first-passage processes* (Cambridge University Press, 2001).
 - [19] N. G. Van Kampen, *Stochastic processes in physics and chemistry*, vol. 1 (Elsevier, 1992).
 - [20] T. A. Fulton and G. J. Dolan, *Phys. Rev. Lett.* **59**, 109 (1987), ISSN 1079-7114.
 - [21] H. Pothier, P. Lafarge, C. Urbina, D. Esteve, and M. H. Devoret, *EPL (Europhysics Letters)* **17**, 249 (1992).
 - [22] M. H. Devoret, D. Estève, H. Grabert, G.-L. Ingold, H. Pothier, and C. Urbina, *Physical review letters* **64**, 1824 (1990).
 - [23] P. Zheng, Ph.D. thesis, UC Berkeley (2016).
 - [24] S. Mukhopadhyay, H. Mahmoodi-Meimand, and K. Roy, in *2004 Symposium on VLSI Circuits. Digest of Technical Papers (IEEE Cat. No. 04CH37525)* (IEEE, 2004), pp. 64–67.
 - [25] T. Cossetto, *Problems in nonequilibrium fluctuations across scales: A path integral approach* (2020), [Online; accessed 19. Feb. 2021], URL <https://orbilu.uni.lu/>

- handle/10993/45484.
- [26] A. Kamenev, *Field Theory of Non-Equilibrium Systems* (Cambridge University Press, Cambridge, England, UK, 2011), ISBN 978-0-52176082-9.
- [27] M. F. Weber and E. Frey, Reports on Progress in Physics **80**, 046601 (2017).
- [28] F. Bouchet and J. Reygner, in *Annales Henri Poincaré* (Springer, 2016), vol. 17, pp. 3499–3532.
- [29] S. Arrhenius, Zeitschrift für physikalische Chemie **4**, 96 (1889).
- [30] H. Eyring, The Journal of Chemical Physics **3**, 107 (1935).
- [31] H. A. Kramers, Physica **7**, 284 (1940).
- [32] N. Berglund, arXiv preprint arXiv:1106.5799 - (2011).
- [33] G. Falasco and M. Esposito, Phys. Rev. E **103**, 042114 (2021).
- [34] K. V. Palem, IEEE Transactions on Computers **54**, 1123 (2005).
- [35] J. Kaiser, R. Faria, K. Y. Camsari, and S. Datta, Frontiers in Computational Neuroscience **14** (2020).
- [36] W. A. Borders, A. Z. Pervaiz, S. Fukami, K. Y. Camsari, H. Ohno, and S. Datta, Nature **573**, 390 (2019).
- [37] C. C. Enz and E. A. Vittoz, John Wiley & Sons Inc **68** (2006).
- [38] Y. Tsidividis and C. McAndrew, *Operation and Modeling of the MOS Transistor* (Oxford Univ. Press, 2011).

Appendix A: Deterministic treatment of the CMOS SRAM cell

In this section we derive the deterministic equations for a CMOS SRAM cell working in the sub-threshold regime. We first consider a single inverter with input voltage v_g and output voltage v , and symmetric powering with voltages $V_{dd} = -V_{ss}$. The current $I_p(v, v_g)$ through the pMOS transistor for given v and v_g is [37]:

$$I_p(v, v_g) = I_0 e^{-V_{th}/V_T} e^{(V_{dd}-v_g)/(nV_T)} (1 - e^{-(V_{dd}-v)/V_T}), \quad (\text{A1})$$

while for the nMOS transistor we have $I_n(v, v_g) = I_p(-v, -v_g)$. From this we can construct the deterministic dynamical equations for the voltages v_1 and v_2 of the CMOS SRAM cell discussed in the main text:

$$\begin{aligned} C \frac{dv_1}{dt} &= I_p(v_1, v_2) - I_n(v_1, v_2) \\ C \frac{dv_2}{dt} &= I_p(v_2, v_1) - I_n(v_2, v_1). \end{aligned} \quad (\text{A2})$$

We first solve for the stationary solution satisfying $dv_1/dt = dv_2/dt = 0$. By symmetry, this solution must satisfy $v_1 = -v_2 = v^*$. Thus, we need to find v^* such that $I_p(v^*, -v^*) = I_n(v^*, -v^*)$. In the following, for simplicity, we consider the case $n = 1$. In that case, the possible solutions are $v_0 = 0$ and, only if $V_{dd} > V_T \log(2)$,

$$v_{\pm} = V_{dd} + V_T \log \left(1/2 \pm \sqrt{1/4 - e^{-2V_{dd}/V_T}} \right). \quad (\text{A3})$$

Note that $v_+ = -v_-$, since actually these are the two solutions in the bistable phase. We now consider $v_1 = v_+ + \delta v_1$ and $v_2 = v_- + \delta v_2$ and expand Eq. (A2) to first order in $\delta v_{1/2}$, finding:

$$\frac{d}{dt} \begin{bmatrix} \delta v_1 \\ \delta v_2 \end{bmatrix} = \frac{I_0 e^{-V_{th}/V_T}}{CV_T} \begin{bmatrix} 2 - e^{2V_{dd}/V_T} & -2 \\ -2 & 2 - e^{2V_{dd}/V_T} \end{bmatrix} \begin{bmatrix} \delta v_1 \\ \delta v_2 \end{bmatrix} \quad (\text{A4})$$

The eigenvalues of the matrix in the previous equation are $-e^{2V_{dd}/V_T}$ and $4 - e^{2V_{dd}/V_T}$, which shows that the solution considered is indeed stable for $V_{dd} > V_T \log(2)$ (a similar analysis shows that the solution v_0 becomes unstable at the same point), and that small departures relax back to it at a rate $\lambda_{eq} \simeq \tau_0^{-1} (v_e/V_T) e^{2V_{dd}/V_T}$, with $\tau_0 = (q_e/I_0) e^{V_{th}/V_T}$.

From the previous solution it can be seen that the stationary current through each transistor is $I_n = I_p = (q_e/\tau_0)(e^{V_{dd}/V_T} - 1)$ for $V_{dd} \leq V_T \log(2)$ (monostability), and $I_n = I_p = q_e/\tau_0$ for $V_{dd} > V_T \log(2)$ (bistability). Thus, the current in the bistable phase is constant and the total entropy production is $\dot{\Sigma} = 2(2V_{dd}q_e/\tau_0)$.

Appendix B: Macroscopic limit and large deviations principle

The conduction channel of a MOS transistor in typical designs has two associated dimensions: its width W and its length L [37, 38]. The capacitance between the gate terminal and the body of the transistor (which is typically the largest one) scales as the area of the channel: $C \propto WL$. Also, the current through the channel for fixed drain-source and gate-source voltages is proportional to the channel width, and inversely proportional to the channel length [37]. Thus, the parameter I_0 used to characterize the I-V curve of the transistor scales as $I_0 \propto W/L$. For the following

discussion we are going to consider a family of devices with fixed channel length, but variable channel width. Thus, we can consider W as a scale parameter, with respect to which both the capacitance C and the current I_0 are proportional. In that case, as considered in the main text, the elementary voltage $v_e = q_e/C$ scales as W^{-1} , while the Poisson rates $\lambda_{\pm}^{n/p}(v_1, v_2)$ associated to the transistors scale as W . Under those conditions, the master equation in the main text can be rewritten as:

$$d_t P(\mathbf{v}, t) = \sum_{\rho} v_e^{-1} [\omega_{\rho}(\mathbf{v} - v_e \Delta_{\rho}, v_e) P(\mathbf{v} - v_e \Delta_{\rho}) - \omega_{\rho}(\mathbf{v}, v_e) P(\mathbf{v})]. \quad (\text{B1})$$

In the previous equation, $\mathbf{v} = (v_1, v_2)^T$ is the state vector, and the index runs over the possible transitions. For example, the values $\rho = 1, \dots, 4$ correspond to the forward transitions of each transistor, while $\rho = -1, \dots, -4$ to the reverse transitions. The vectors Δ_{ρ} encode the change in voltage associated to each transition. The scaled rates $\omega_{\rho}(\mathbf{v}, v_e)$ are related to the original Poisson rates $\omega_{\rho}(\mathbf{v}, v_e)$ by $\lambda_{\rho}(\mathbf{v}, v_e) = v_e^{-1} \omega_{\rho}(\mathbf{v}, v_e)$. Thus, the scaling of the rates with respect to W (or equivalently, with respect to v_e), is taken into account in the factor v_e^{-1} , in such a way that the limit $\lim_{v_e \rightarrow 0} \omega_{\rho}(\mathbf{v}, v_e)$ is well defined (the limit $v_e \rightarrow 0$ here and below must be interpreted as $v_e/V_T \rightarrow 0$ and $v_e/V_{\text{dd}} \rightarrow 0$ for fixed V_T and V_{dd}). Note that the explicit dependence of the rates in the elementary voltage v_e stems from the charging effects discussed in the main text.

Under these conditions, the solution of the master equation in Eq. (B1) satisfies a large deviations principle in the macroscopic limit $v_e \rightarrow 0$. In order to see this, we introduce the large deviations ansatz $P(\mathbf{v}, t) \asymp \exp(-(f(\mathbf{v}, t) + o(v_e))/v_e)$ into Eq. (B1), and only keep the dominant terms in $v_e \rightarrow 0$. We note that in that limit $P(\mathbf{v} - v_e \Delta_{\rho}, t) \asymp P(\mathbf{v}, t) \exp((\Delta_{\rho})_i \partial_{v_i} f(\mathbf{v}, t))$. Therefore, the master equation in Eq. (B1) reduces to the following dynamical equation for the rate function:

$$d_t f(\mathbf{v}, t) = \sum_{\rho} \omega_{\rho}(\mathbf{v}, 0) \left[1 - e^{(\Delta_{\rho})_i \partial_{v_i} f(\mathbf{v}, t)} \right]. \quad (\text{B2})$$

It is worth noting that for general jump processes with scaling properties as the ones satisfied by Eq. (B1), the validity of the large deviation principle can be formally proven [25].

For the particular circuit under consideration, we can see from the previous equation that the steady state rate function $f(v_1, v_2)$ should satisfy

$$0 = (e^{\partial_{v_1} f} - 1) a(v_1, v_2) + (e^{-\partial_{v_1} f} - 1) b(v_1, v_2) + (e^{\partial_{v_2} f} - 1) a(v_2, v_1) + (e^{-\partial_{v_2} f} - 1) b(v_2, v_1), \quad (\text{B3})$$

as presented in the main text, where the functions $a(v_1, v_2)$ and $b(v_1, v_2)$ were defined as the appropriate combination of the scaled transition rates. The previous equation cannot be solved exactly. However, it can be employed to solve for reduced rate functions derived from $f(v_1, v_2)$, exploiting the symmetry of the problem and the contraction principle of large deviations theory. We begin by changing variables to $x = (v_1 - v_2)/2$ and $y = (v_1 + v_2)/2$. Then, $\partial_{1/2} f = (\pm \partial_x f + \partial_y f)/2$. Defining $\alpha = e^{\partial_x f/2}$ and $\beta = e^{\partial_y f/2}$, the previous equation becomes:

$$0 = (\alpha\beta - 1) a(x, y) + (\alpha^{-1}\beta^{-1} - 1) b(x, y) + (\beta/\alpha - 1) a(-x, y) + (\alpha/\beta - 1) b(-x, y), \quad (\text{B4})$$

where the change of variables of the functions $a(x, y)$ and $b(x, y)$ is implicit. Now, we are interested in computing the partial distributions $P(x)$ and $Q(y)$ for the variables x and y . The contraction principle states that if the full distribution $P_{\text{ss}}(x, y)$ satisfies a large deviation principle with rate function $f(x, y)$, then the partial distribution $P(x) = \sum_y P_{\text{ss}}(x, y)$ also satisfies a large deviation principle with rate function $g(x) = \inf_y f(x, y)$ [14]. Then, assuming that f is sufficiently regular and that $\inf_y f(x, y) = \min_y f(x, y)$, we have $g(x) = f(x, y_{\min}|x)$, where $y_{\min}|x$ is a minimum of $f(x, \cdot)$ and therefore satisfies $\partial_y f(x, y_{\min}|x) = 0$. Thus, $y_{\min}|x$ is the most probable value of y for a given value of x . As discussed in the main text, the symmetry of the steady state is such that $y_{\min}|x = 0$ for all x . Thus, evaluating the previous equation at $y = 0$, since $\beta|_{y=0} = 1$, we obtain:

$$\alpha|_{y=0} = e^{\partial_x f|_{y=0}/2} = e^{d_x g(x)/2} = \frac{a(-x, 0) + b(x, 0)}{a(x, 0) + b(-x, 0)}, \quad (\text{B5})$$

from where we easily obtain the expression for $d_x g(x)$ given in the main text. Note that from the previous expression is evident that $d_x g(x)$ is an odd function, and therefore $g(x)$ is even. The rate function $h(y)$ for the partial distribution $Q(y)$ can also be obtained. It is given by $h(y) = f(x_{\min}|y, y)$, where $x_{\min}|y$ is a minimum of $f(\cdot, y)$. In this case a symmetry argument is lacking, and to proceed we must neglect correlations between x and y . Then, x_{\min} is considered to be independent of y , and thus it can be computed as the minimum of $g(x)$. In the bistable phase there are actually

two equivalent values of x_{\min} , that lead to the same function $h(y)$. Thus, evaluating Eq. (B4) at $x = x_{\min}$, since $\alpha|_{x=x_{\min}} = 1$, we obtain:

$$\beta|_{x=x_{\min}} = e^{\partial_y f|_{x=x_{\min}}/2} = e^{d_y h(y)/2} = \frac{b(x_{\min}, y) + b(-x_{\min}, y)}{a(x_{\min}, y) + a(-x_{\min}, y)}. \quad (\text{B6})$$

For the Poisson rates corresponding to MOS transistors in subthreshold operation, that enter into the definition of the functions $a(x, y)$ and $b(x, y)$, the integration of $d_x g(x)$ can be performed exactly, leading to the compact expression given in the main text. This is not the case for $d_y h(y)$. However, it is possible to obtain the leading behaviour of $h(y)$ around $y = 0$, which is given by:

$$h(y) = \frac{2}{n} \frac{(n-1)(1 + e^{2(1+1/n)x_{\min}/V_T}) + e^{(V_{\text{dd}}+x_{\min})/V_T} + e^{(V_{\text{dd}}+x_{\min}(1+2/n))/V_T}}{1 + e^{2(1+1/n)x_{\min}/V_T} + e^{(V_{\text{dd}}+x_{\min})/V_T} + e^{(V_{\text{dd}}+x_{\min}(1+2/n))/V_T}} y^2/V_T + \mathcal{O}(y^4). \quad (\text{B7})$$

Finally, we note that the most probable values according to the large deviations solution ($x = x_{\min}$ and $y = 0$, which correspond to $v_1 = -v_2 = x_{\min}$) match the deterministic solutions obtained in the previous section.

Appendix C: Validity of the separability assumption

In this section we discuss the accuracy of the separability assumption used above to derive equation Eq. (B6) for the rate function $h(y)$, and also involved in the reconstruction of the full probability distribution in Eq. (8) in the main text. We stress the fact that the expression for the rate function $g(x)$ is exact and independent of such assumption.

We begin by considering the regime of small fluctuations around the deterministic fixed points. We first note that for small perturbations around the fixed points, the deterministic dynamics completely decouples the variables $x = (v_1 - v_2)/2$ and $y = (v_1 + v_2)/2$, as can be seen from Eq. (A4). However, fluctuations might still induce correlations. To see that this is not the case for typical fluctuations, we will compute the Gaussian fluctuations around the fixed points by expanding Eq. (B2). Thus, if \mathbf{v}^* are the fixed point voltages, the rate function $f(\mathbf{v})$ can be expanded as:

$$f(\mathbf{v}) = \frac{1}{2}(\mathbf{v} - \mathbf{v}^*)^T C(\mathbf{v} - \mathbf{v}^*) + \mathcal{O}(|\mathbf{v} - \mathbf{v}^*|^3) \quad (\text{C1})$$

in terms of the matrix $\{C\}_{ij} = d_{v_i v_j}^2 f(\mathbf{v}^*)$. Then, the Gaussian covariance matrix is given by $v_e C^{-1}$. Accordingly, the covariance matrix for the variables x and y is $v_e M C^{-1} M^T$, with $M = \begin{pmatrix} 1/2 & -1/2 \\ 1/2 & 1/2 \end{pmatrix}$. We will show now that $M C^{-1} M^T$ is a diagonal matrix, and therefore the variables x and y are uncorrelated to the Gaussian level. Expanding Eq. (B2) to second order in $\mathbf{v} - \mathbf{v}^*$ we obtain:

$$0 = C^{-1} A + A^T C^{-1} + B, \quad (\text{C2})$$

where the matrices A and B are given by

$$\{A\}_{ij} = \sum_{\rho} \partial_{v_i} \omega_{\rho}(\mathbf{v}^*, 0) (\Delta_{\rho})_j \quad \text{and} \quad \{B\}_{ij} = \sum_{\rho} \omega_{\rho}(\mathbf{v}^*, 0) (\Delta_{\rho})_i (\Delta_{\rho})_j. \quad (\text{C3})$$

In this particular model one can see that $B = b\mathbb{1}$ is always proportional to the identity matrix. In particular for $n = 1$ we have $b = 4(v_e/\tau_0)e^{V_{\text{dd}}/V_T} \sinh(V_{\text{dd}}/V_T)$. Also, the matrix A is the one appearing in Eq. (A4) (in the bistable phase, i.e., for $V_{\text{dd}} \geq V_T \log(2)$). The matrix $C'^{-1} = M C^{-1} M^T$ satisfies an equation analogous to Eq. (C2), but in terms of the transformed matrices $A' = (M^T)^{-1} A M^T$ and $B' = M B M^T$. Solving that equation one can easily show that the matrix C^{-1} is indeed diagonal.

The difference between the exact steady state distribution $P_{\text{ss}}^{\text{ex}}(\mathbf{v})$ and the approximated reconstruction $P_{\text{ss}}(\mathbf{v})$ given by Eq. (8) in the main text can be quantified by the Hellinger distance $0 \leq H \leq 1$, which is computed as $H^2 = 1 - \sum_{\mathbf{v}} \sqrt{P_{\text{ss}}^{\text{ex}}(\mathbf{v}) P_{\text{ss}}(\mathbf{v})}$. Since the separability assumption holds for Gaussian fluctuations, and they are the dominant ones in the macroscopic limit $v_e/V_T \rightarrow 0$, it follows that $H \rightarrow 0$ in the same limit. For finite values of v_e , the accuracy of the reconstruction of Eq. (8) can be tested numerically. As an example, in Figure 5 we compare the exact steady state distribution $P_{\text{ss}}^{\text{ex}}(\mathbf{v})$ with the reconstruction $P_{\text{ss}}(\mathbf{v})$ for $V_{\text{dd}}/V_T = 1.2$ and $v_e/V_T = 0.1$, and we also show the difference $P_{\text{ss}}^{\text{ex}}(\mathbf{v}) - P_{\text{ss}}(\mathbf{v})$. In this case, the Hellinger distance is $H = 1.94 \times 10^{-2} \ll 1$.

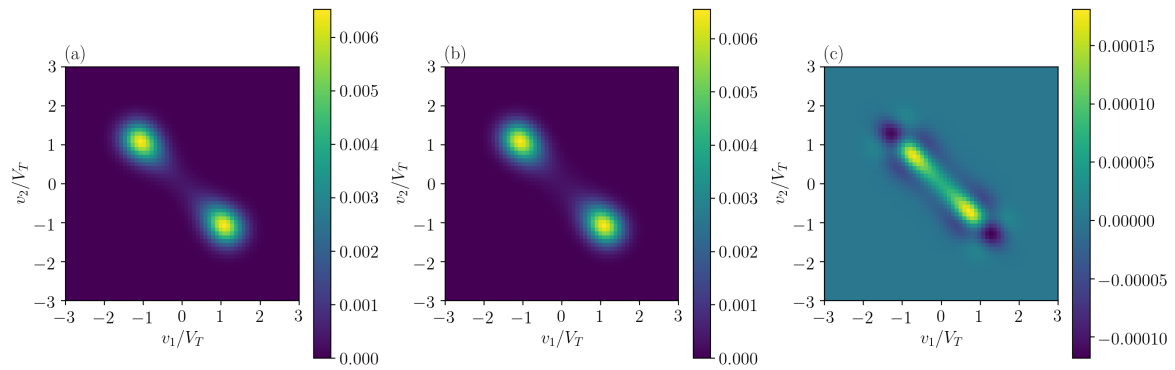


FIG. 5: (a) Exact steady state distribution $P_{ss}^{ex}(\mathbf{v})$ for $V_{dd}/V_T = 1.2$ and $v_e/V_T = 0.1$. (b) Reconstruction $P_{ss}(\mathbf{v})$ based on Eq. (8) and the analytical expressions for the rate functions $g(x)$ and $h(y)$, for the same parameters. (c) Difference $P_{ss}^{ex}(\mathbf{v}) - P_{ss}(\mathbf{v})$.