

Fast Statistical Analysis of Rare Failure Events for Memory Circuits in High-Dimensional Variation Space

Shupeng Sun and Xin Li

Electrical & Computer Engineering Department, Carnegie Mellon University
5000 Forbes Avenue, Pittsburgh, PA 15213 USA
{shupengs, xinli}@ece.cmu.edu

ABSTRACT - Accurately estimating the rare failure rates for nanoscale memory circuits is a challenging task, especially when the variation space is high-dimensional. In this paper, we summarize two novel techniques to address this technical challenge. First, we describe a subset simulation (SUS) technique to estimate the rare failure rates for continuous performance metrics. The key idea of SUS is to express the rare failure probability of a given circuit as the product of several large conditional probabilities by introducing a number of intermediate failure events. These conditional probabilities can be efficiently estimated with a set of Markov chain Monte Carlo samples generated by a modified Metropolis algorithm. Second, to efficiently estimate the rare failure rates for discrete performance metrics, scaled-sigma sampling (SSS) can be used. SSS aims to generate random samples from a distorted probability distribution for which the standard deviation (i.e., sigma) is scaled up. Next, the failure rate is accurately estimated from these scaled random samples by using an analytical model derived from the theorem of “soft maximum”. Our experimental results of several nanoscale circuit examples demonstrate that SUS and SSS achieve significantly improved accuracy over other traditional techniques when the dimensionality of the variation space is more than a few hundred.

I. Introduction

As integrated circuit (IC) technology advances, the ever increasing process variation has become a growing concern [1]. A complex IC, containing numerous memory components, is required to meet the design specification not only at the nominal process corner, but also under large-scale process variations. To achieve sufficiently high yield, the failure rate of each individual memory component must be extremely small. For instance, the failure rate of an SRAM bit-cell must be less than 10^{-8} ~ 10^{-6} for a typical SRAM design [2]-[3]. Due to this reason, efficiently analyzing the rare failure event for the individual memory component becomes an important task for the IC design community.

The simple way to estimate the failure probability is to apply the well-known crude Monte Carlo (CMC) technique [20]. CMC directly draws random samples from the probability density function (PDF) that models device-level variations, and performs a transistor-level simulation to evaluate the performance value for each sample. When CMC is applied to estimate an extremely small failure rate (e.g., 10^{-8} ~ 10^{-6}), most random samples do not fall into the failure region. Hence, a large number of (e.g., 10^7 ~ 10^9) samples are needed to accurately estimate the small failure probability, which implies that CMC can be extremely expensive for our application of rare failure rate estimation.

To improve the sampling efficiency, importance sampling (IS) methods have been proposed in the literature [4], [7], [10], [12]-[13]. Instead of sampling the original PDF, IS samples a distorted PDF to get more samples in the important failure region. The efficiency achieved by IS highly depends

on the choice of the distorted PDF. The traditional IS methods apply several heuristics to construct a distorted PDF that can capture the most important failure region in the variation space. Such a goal, though easy to achieve in a low-dimensional variation space, is extremely difficult to fulfill when a large number of random variables are used to model process variations.

Another approach to improving the sampling efficiency, referred to as statistical blockade, has recently been proposed [9]. This approach first builds a classifier with a number of transistor-level simulations, and then draws random samples from the original PDF. Unlike CMC where all the samples are evaluated by transistor-level simulations, statistical blockade only simulates the samples that are likely to fall into the failure region or close to the failure boundary based on the classifier. The efficiency achieved by this approach highly depends on the accuracy of the classifier. If the variation space is high-dimensional, a large number of transistor-level simulations are needed to build an accurate classifier, which makes the statistical blockade method quickly intractable.

In addition to the aforementioned statistical methods, several deterministic approaches have also been proposed to efficiently estimate the rare failure probability [6], [14]. These methods first find the failure boundary, and then calculate the failure probability by integrating the PDF over the failure region in the variation space. Though efficient in a low-dimensional variation space, it is often computationally expensive to accurately determine the failure boundary in a high-dimensional space especially if the boundary has a complicated shape (e.g., non-convex or even discontinuous).

Most of these traditional methods [4]-[14] have been successfully applied to SRAM bit-cells to estimate their rare failure rates where only a small number of (e.g., 6~20) independent random variables are used to model process variations and, hence, the corresponding variation space is low-dimensional. It has been demonstrated in the literature that estimating the rare failure probability in a high-dimensional space (e.g., hundreds of independent random variables to model the device-level variations for SRAM) becomes increasingly important [18]. Unfortunately, such a high-dimensional problem cannot be efficiently handled by most traditional methods. It, in turn, poses an immediate need of developing a new CAD tool to accurately capture the rare failure events in a high-dimensional variation space with low computational cost.

To address this technical challenge, we first describe a novel subset simulation (SUS) technique. The key idea of SUS, borrowed from the statistics community [15]-[17], is to express the rare failure probability as the product of several large conditional probabilities by introducing a number of intermediate failure events. As such, the original problem of rare failure probability estimation is cast to an equivalent

problem of estimating a sequence of conditional probabilities via multiple phases. Since these conditional probabilities are relatively large, they are substantially easier to estimate than the original rare failure rate.

When implementing the SUS method, it is difficult, if not impossible, to directly draw random samples from the conditional PDFs and estimate the conditional probabilities, since these conditional PDFs are unknown in advance. To address this issue, a modified Metropolis (MM) algorithm is adopted from the literature [15] to generate random samples by constructing a number of Markov chains. The conditional probabilities of interest are then estimated from these random samples. Unlike most traditional techniques [4]-[14] that suffer from the dimensionality issue, SUS can be efficiently applied to high-dimensional problems, which will be demonstrated by the experimental results in Section II.

To define the intermediate failure events required by SUS, the performance of interest (PoI) must be continuous. In other words, SUS can only analyze a continuous PoI. For many rare failure events, however, PoIs are discrete (e.g., the output of a voltage-mode sense amplifier). Realizing this limitation, we further describe a scaled-sigma sampling (SSS) approach to efficiently estimate the rare failure rates for discrete PoIs in a high-dimensional space. SSS is particularly developed to address the following two fundamental questions: (i) how to efficiently draw random samples from the rare failure region, and (ii) how to estimate the rare failure rate based on these random samples. Unlike CMC that directly samples the variation space and therefore only few samples fall into the failure region, SSS draws random samples from a distorted PDF for which the standard deviation (i.e., sigma) is scaled up. Conceptually, it is equivalent to increasing the magnitude of process variations. As a result, a large number of samples can now fall into the failure region. Once the distorted random samples are generated, an analytical model derived from the theorem of “soft maximum” is optimally fitted by applying maximum likelihood estimation (MLE). Next, the failure rate can be efficiently estimated from the fitted model.

The remainder of this paper is organized as follows. In Section II, we will summarize the SUS approach, and then the SSS approach will be presented in Section III. Finally, we conclude in Section IV.

II. Subset Simulation

Suppose that the vector

$$\mathbf{x} = [x_1 \ x_2 \ \cdots \ x_M]^T \quad (1)$$

is an M -dimensional random variable modeling device-level process variations and its joint PDF is $f(\mathbf{x})$. The failure rate of a circuit can be mathematically represented as

$$P_F = \Pr(\mathbf{x} \in \Omega) = \int_{\mathbf{x} \in \Omega} f(\mathbf{x}) \cdot d\mathbf{x}, \quad (2)$$

where Ω denotes the failure region, i.e., the subset of the variation space where the performance of interest does not meet the specification.

Instead of directly estimating the rare failure probability P_F , SUS expresses P_F as the product of several large conditional probabilities by introducing a number of intermediate failure events in the variation space. Without loss of generality, we define K intermediate failure events $\{\Omega_k; k = 1, 2, \dots, K\}$ as

$$\Omega_1 \supset \Omega_2 \supset \cdots \supset \Omega_{K-1} \supset \Omega_K = \Omega. \quad (3)$$

Based on (3), we can express P_F in (2) as

$$P_F = \Pr(\mathbf{x} \in \Omega) = \Pr(\mathbf{x} \in \Omega_K, \mathbf{x} \in \Omega_{K-1}). \quad (4)$$

Eq. (4) can be re-written as

$$P_F = \Pr(\mathbf{x} \in \Omega_K | \mathbf{x} \in \Omega_{K-1}) \cdot \Pr(\mathbf{x} \in \Omega_{K-1}). \quad (5)$$

Similarly, we can express $\Pr(\mathbf{x} \in \Omega_{K-1})$ as

$$\Pr(\mathbf{x} \in \Omega_{K-1}) = \Pr(\mathbf{x} \in \Omega_{K-1} | \mathbf{x} \in \Omega_{K-2}) \cdot \Pr(\mathbf{x} \in \Omega_{K-2}). \quad (6)$$

From (3), (5) and (6), we can easily derive

$$P_F = \Pr(\mathbf{x} \in \Omega_1) \cdot \prod_{k=2}^K \Pr(\mathbf{x} \in \Omega_k | \mathbf{x} \in \Omega_{k-1}) = \prod_{k=1}^K P_k, \quad (7)$$

where

$$P_1 = \Pr(\mathbf{x} \in \Omega_1) \quad (8)$$

$$P_k = \Pr(\mathbf{x} \in \Omega_k | \mathbf{x} \in \Omega_{k-1}) \quad (k = 2, 3, \dots, K). \quad (9)$$

If $\{\Omega_k; k = 1, 2, \dots, K\}$ are properly chosen, all the probabilities $\{P_k; k = 1, 2, \dots, K\}$ are large and can be efficiently estimated. Once $\{P_k; k = 1, 2, \dots, K\}$ are known, the rare failure probability P_F can be easily calculated by (7).

Note that the failure events $\{\Omega_k; k = 1, 2, \dots, K\}$ are extremely difficult to specify in a high-dimensional variation space. For this reason, we do not directly define $\{\Omega_k; k = 1, 2, \dots, K\}$ in the variation space. Instead, we utilize their corresponding subsets $\{F_k; k = 1, 2, \dots, K\}$ in the performance space

$$F_k = \{y(\mathbf{x}); \mathbf{x} \in \Omega_k\} \quad (k = 1, 2, \dots, K), \quad (10)$$

where $y(\mathbf{x})$ denotes the performance of interest (PoI) as a function of \mathbf{x} . Since $y(\mathbf{x})$ is typically a scalar, $\{F_k; k = 1, 2, \dots, K\}$ are just one-dimensional subsets of \mathbb{R} and, therefore, easy to specify. Once $\{F_k; k = 1, 2, \dots, K\}$ are determined, $\{\Omega_k; k = 1, 2, \dots, K\}$ are implicitly known. For instance, to know whether a given \mathbf{x} belongs to Ω_k , we first run a transistor-level simulation to evaluate $y(\mathbf{x})$. If $y(\mathbf{x})$ belongs to F_k , \mathbf{x} is inside Ω_k . Otherwise, \mathbf{x} is outside Ω_k .

In what follows, we will use a simple 2-D example to intuitively illustrate the basic flow of SUS. Fig. 1 shows this 2-D example where two random variables $\mathbf{x} = [x_1 \ x_2]$ are used to model the device-level process variations, and Ω_1 and Ω_2 denote the first two subsets in (3). Note that Ω_1 and Ω_2 are depicted for illustration purposes in this example. In practice, we do not need to explicitly know Ω_1 and Ω_2 , as previously explained.

Our objective is to estimate the probabilities $\{P_k; k = 1, 2, \dots, K\}$ via multiple phases. Starting from the 1st phase, we simply draw L_1 independent random samples $\{\mathbf{x}^{(1,l)}; l = 1, 2, \dots, L_1\}$ from the PDF $f(\mathbf{x})$ to estimate P_1 . Here, the superscript “1” of the symbol $\mathbf{x}^{(1,l)}$ refers to the 1st phase. Among these L_1 samples, we identify a subset of samples $\{\mathbf{x}_F^{(1,t)}; t = 1, 2, \dots, T_1\}$ that fall into Ω_1 , where T_1 denotes the total number of samples in this subset. As shown in Fig. 1 (a), the red points represent the samples that belong to Ω_1 and the green points represent the samples that are out of Ω_1 . In this case, P_1 can be estimated as

$$P_1^{SUS} = \frac{1}{L_1} \cdot \sum_{t=1}^{T_1} I_{\Omega_1}[\mathbf{x}^{(1,t)}] = \frac{T_1}{L_1}, \quad (11)$$

where P_1^{SUS} denotes the estimated value of P_1 , and $I_{\Omega_1}(\mathbf{x})$

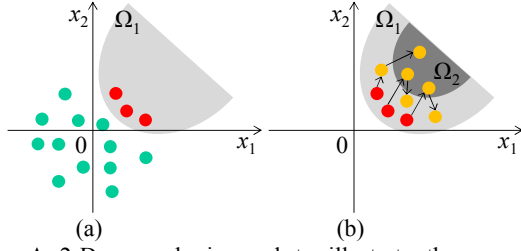


Fig. 1. A 2-D example is used to illustrate the procedure of probability estimation via multiple phases by using SUS: (a) generating MC samples and estimating P_1 in the 1st phase, and (b) generating MCMC samples and estimating P_2 in the 2nd phase.

represents the indicator function

$$I_{\Omega_1}(\mathbf{x}) = \begin{cases} 1 & \mathbf{x} \in \Omega_1 \\ 0 & \mathbf{x} \notin \Omega_1 \end{cases}. \quad (12)$$

If P_1 is large, it can be accurately estimated with a small number of random samples (e.g., L_1 is around $10^2 \sim 10^3$).

Next, in the 2nd phase, we need to estimate the conditional probability $P_2 = \Pr(\mathbf{x} \in \Omega_2 | \mathbf{x} \in \Omega_1)$. Towards this goal, one simple idea is to directly draw random samples from the conditional PDF $f(\mathbf{x} | \mathbf{x} \in \Omega_1)$ and then compute the mean of the indicator function $I_{\Omega_2}(\mathbf{x})$

$$I_{\Omega_2}(\mathbf{x}) = \begin{cases} 1 & \mathbf{x} \in \Omega_2 \\ 0 & \mathbf{x} \notin \Omega_2 \end{cases}. \quad (13)$$

This approach, however, is practically infeasible since $f(\mathbf{x} | \mathbf{x} \in \Omega_1)$ is unknown in advance. To address this issue, we apply a modified Metropolis (MM) algorithm [15] to generate a set of random samples that follow the conditional PDF $f(\mathbf{x} | \mathbf{x} \in \Omega_1)$.

MM is a Markov chain Monte Carlo (MCMC) technique [20]. To clearly describe the sampling procedure of MM, we define the symbol $\mathbf{x}^{(2,t,1)} = \mathbf{x}_F^{(1,t)}$, where $t \in \{1, 2, \dots, T_1\}$. The superscripts “2” and “1” of $\mathbf{x}^{(2,t,1)}$ refer to the 2nd phase and the 1st sample of the Markov chain respectively. Starting from each of the samples $\{\mathbf{x}^{(2,t,1)}; t = 1, 2, \dots, T_1\}$, MM generates a sequence of samples that form a Markov chain $\{\mathbf{x}^{(2,t,l)}; l = 1, 2, \dots, L_2\}$, where L_2 denotes the length of the Markov chain in the 2nd phase. T_1 Markov chains are independently generated by MM, and in total we obtain $T_1 \cdot L_2$ MCMC samples: $\{\mathbf{x}^{(2,t,l)}; t = 1, 2, \dots, T_1, l = 1, 2, \dots, L_2\}$. Due to the page limit, more details about how to generate a Markov chain from an initial sample via MM are omitted in this paper, but can be found in the literature [15].

Fig. 1 (b) shows the sampling results in the 2nd phase for our 2-D example. In Fig. 1 (b), the red points represent the initial samples $\{\mathbf{x}^{(2,t,1)}; t = 1, 2, \dots, T_1\}$ of the Markov chains and they are obtained from the 1st phase. The yellow points represent the MCMC samples created via the MM algorithm in the 2nd phase. It has been proved in [15] that all these MCMC samples $\{\mathbf{x}^{(2,t,l)}; t = 1, 2, \dots, T_1, l = 1, 2, \dots, L_2\}$ in Fig. 1 (b) approximately follow $f(\mathbf{x} | \mathbf{x} \in \Omega_1)$. In other words, we have successfully generated a number of random samples that follow our desired distribution for the 2nd phase.

Among all the MCMC samples $\{\mathbf{x}^{(2,t,l)}; t = 1, 2, \dots, T_1, l = 1, 2, \dots, L_2\}$, we further identify a subset of samples $\{\mathbf{x}_F^{(2,t,l)}; t = 1, 2, \dots, T_2\}$ that fall into Ω_2 , where T_2 denotes the total number of the samples in this subset. The conditional probability P_2 can be estimated as

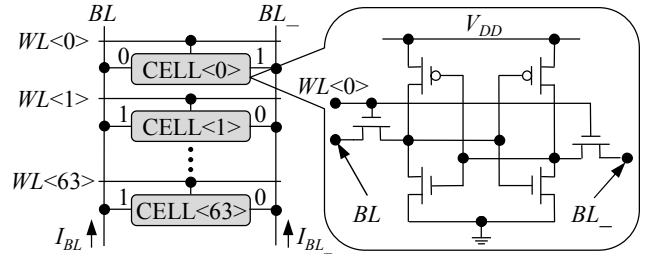


Fig. 2. The simplified schematic is shown for an SRAM column consisting of 64 bit-cells designed in a 45nm CMOS process.

$$P_2^{SUS} = \frac{1}{T_1 \cdot L_2} \cdot \sum_{t=1}^{T_1} \sum_{l=1}^{L_2} I_{\Omega_2}[\mathbf{x}^{(2,t,l)}] = \frac{T_2}{T_1 \cdot L_2}, \quad (14)$$

where P_2^{SUS} denotes the estimated value of P_2 .

By following the aforementioned idea, we can estimate all the probabilities $\{P_k; k = 1, 2, \dots, K\}$. Once the values of $\{P_k; k = 1, 2, \dots, K\}$ are estimated, the rare failure rate P_F is calculated by

$$P_F^{SUS} = \prod_{k=1}^K P_k^{SUS}, \quad (15)$$

where P_F^{SUS} represents the estimated value of P_F by using SUS. If we have more than two random variables, estimating the probabilities $\{P_k; k = 1, 2, \dots, K\}$ can be pursued in a similar way [19].

To efficiently apply SUS, the subsets $\{F_k; k = 1, 2, \dots, K\}$ must be carefully chosen. In addition, to make SUS of practical utility, the statistical property of the SUS estimator must be accurately evaluated so that the estimation error of SUS can be calculated. Due to the page limit, all these implementation issues are not discussed in this paper, but can be found in a recent publication [19].

To demonstrate the efficacy of SUS, we consider an SRAM column example designed in a 45nm CMOS process, as shown in Fig. 2. In this example, our PoI is the read current I_{READ} , which is defined as the difference between the bit-line currents I_{BL} and I_{BL_-} (i.e., $I_{READ} = I_{BL} - I_{BL_-}$) when we start to read CELL<0>. If I_{READ} is greater than a pre-defined specification, we consider the SRAM circuit as “PASS”. For process variation modeling, the local V_{TH} mismatch of each transistor is considered as an independent Normal random variable. In total, we have 384 independent random variables (i.e., 64 bit-cells \times 6 transistors per bit-cell = 384).

We first run CMC with 10^9 random samples, and the estimated failure rate is 1.1×10^{-6} , which is considered as the “golden” failure rate in this example. Next, we compare SUS with the traditional importance sampling technique: MNIS [10], where 2000 simulations are used to construct the distorted PDF. We repeatedly run MNIS and SUS for 100 times with 6000 transistor-level simulations in each run. Fig. 3 shows the 100 estimated 95% CIs for each method, where each blue bar represents the CI of a single run, and the red line represents the “golden” failure rate.

In this example, only a single CI estimated from 100 repeated runs by MNIS can cover the “golden” failure rate, implying that MNIS fails to estimate the CIs accurately. This is an important limitation of MNIS, and generally most of the importance sampling techniques, since the user cannot reliably know the actual “confidence” of the estimator in practice. For the SUS approach, however, there are 95 CIs out

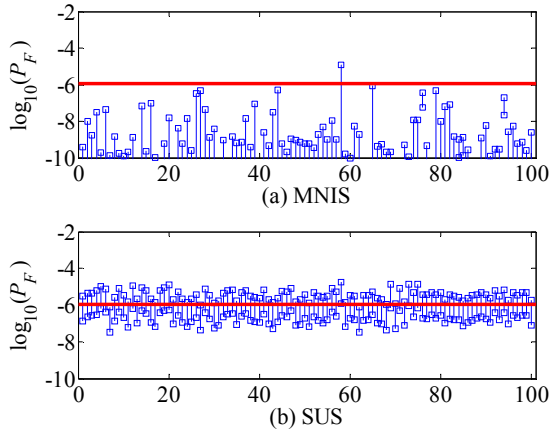


Fig. 3. The 95% confidence intervals (blue bars) of the SRAM read current example are estimated from 100 repeated runs with 6000 transistor-level simulations in each run for: (a) MNIS and (b) SUS. The red line represents the “golden” failure rate.

of 100 runs that cover the “golden” failure rate. More importantly, the CIs estimated by SUS are relatively tight, which implies that SUS achieves substantially better accuracy than the traditional MNIS approach in this example.

Before ending this section, we would like to emphasize that to define the subsets $\{F_k; k = 1, 2, \dots, K\}$ required by SUS, PoI must be continuous. Realizing this limitation, we further describe a scaled-sigma sampling (SSS) approach to efficiently estimate the rare failure rates for discrete PoIs in a high-dimensional space, which will be presented in the next section.

III. Scaled-Sigma Sampling

Unlike the traditional importance sampling methods that must explicitly identify the high-probability failure region, SSS takes a completely different strategy to address the following questions: (i) how to efficiently draw random samples from the high-probability failure region, and (ii) how to estimate the failure rate based on these random samples. In what follows, we will derive the mathematical formulation of SSS and highlight its novelties.

In a process design kit, the random variables $\{x_m; m = 1, 2, \dots, M\}$ in (1) are typically modeled as a jointly Normal distribution [4]-[14]. Without loss of generality, we further assume that $\{x_m; m = 1, 2, \dots, M\}$ are mutually independent and standard Normal (i.e., with zero mean and unit variance)

$$f(\mathbf{x}) = \prod_{m=1}^M \left[\frac{1}{\sqrt{2\pi}} \cdot \exp\left(-\frac{x_m^2}{2}\right) \right] = \frac{\exp(-\|\mathbf{x}\|_2^2/2)}{(\sqrt{2\pi})^M}, \quad (16)$$

where $\|\bullet\|_2$ denotes the L₂-norm of a vector. Any correlated random variables that are jointly Normal can be transformed to the independent random variables $\{x_m; m = 1, 2, \dots, M\}$ by principal component analysis [20].

For the application of rare failure rate estimation, a failure event often occurs at the tail of the PDF $f(\mathbf{x})$. Given (16), it implies that the failure region Ω is far away from the origin $\mathbf{x} = \mathbf{0}$, as shown in Fig. 4 (a). Since the failure rate is extremely small, the traditional CMC analysis cannot efficiently draw random samples from the failure region. Namely, many samples cannot reach the tail of $f(\mathbf{x})$.

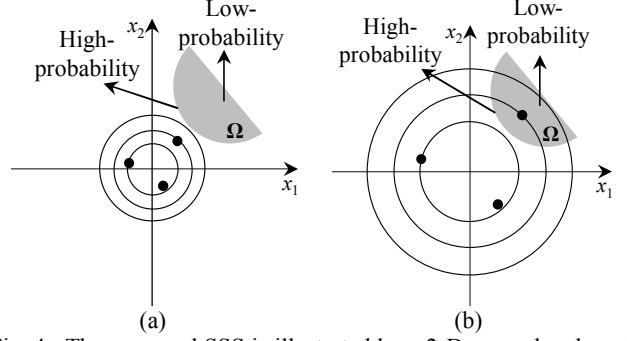


Fig. 4. The proposed SSS is illustrated by a 2-D example where the grey area Ω denotes the failure region and the circles represent the contour lines of the PDF. (a) Rare failure events occur at the tail of the original PDF $f(\mathbf{x})$ and the failure region is far away from the origin $\mathbf{x} = \mathbf{0}$. (b) The scaled PDF $g(\mathbf{x})$ widely spreads over a large region and the scaled samples are likely to reach the far-away failure region.

To address the aforementioned sampling issue, SSS applies a simple idea. Given $f(\mathbf{x})$ in (16), we scale up the standard deviation of \mathbf{x} by a *scaling factor* s ($s > 1$), yielding the following distribution

$$g(\mathbf{x}) = \prod_{m=1}^M \left[\frac{\exp(-x_m^2/2s^2)}{\sqrt{2\pi}s} \right] = \frac{\exp(-\|\mathbf{x}\|_2^2/2s^2)}{(\sqrt{2\pi} \cdot s)^M}. \quad (17)$$

Once the standard deviation of \mathbf{x} is increased by a factor of s , we conceptually increase the magnitude of process variations. Hence, the PDF $g(\mathbf{x})$ widely spreads over a large region and the probability for a random sample to reach the far-away failure region increases, as shown in Fig. 4 (b).

It is important to note that the mean of the scaled PDF $g(\mathbf{x})$ remains $\mathbf{0}$, which is identical to the mean of the original PDF $f(\mathbf{x})$. Hence, for a given sampling location \mathbf{x} , the likelihood defined by the scaled PDF $g(\mathbf{x})$ remains inversely proportional to the length of the vector \mathbf{x} (i.e., $\|\mathbf{x}\|_2$). Namely, it is more (or less) likely to reach the sampling location \mathbf{x} , if the distance between the location \mathbf{x} and the origin $\mathbf{0}$ is smaller (or larger). It, in turn, implies that the high-probability failure region associated with the original PDF $f(\mathbf{x})$ remains the high-probability failure region after the PDF is scaled to $g(\mathbf{x})$, as shown in Fig. 4 (a) and (b). Scaling the PDF from $f(\mathbf{x})$ to $g(\mathbf{x})$ does not change the location of the high-probability failure region; instead, it only makes the failure region easy to sample.

Once the scaled random samples are drawn from $g(\mathbf{x})$ in (17), we need to further estimate the failure rate P_F defined in (2). To this end, one straightforward way is to apply the importance sampling method [20]. Such a simple approach, however, has been proved to be intractable when the dimensionality (i.e., M) of the variation space is high [18]. Namely, it does not fit the need of high-dimensional failure rate estimation in this paper.

Instead of relying on the theory of importance sampling, SSS attempts to estimate the failure rate P_F from a completely different avenue. We first take a look at the “scaled” failure rate P_G corresponding to $g(\mathbf{x})$

$$P_G = \int_{\mathbf{x} \in \Omega} g(\mathbf{x}) \cdot d\mathbf{x} = \int_{-\infty}^{+\infty} I_{\Omega}(\mathbf{x}) \cdot g(\mathbf{x}) \cdot d\mathbf{x}, \quad (18)$$

where $I_{\Omega}(\mathbf{x})$ represents the indicator function

$$I_{\Omega}(\mathbf{x}) = \begin{cases} 1 & \mathbf{x} \in \Omega \\ 0 & \mathbf{x} \notin \Omega \end{cases}. \quad (19)$$

Our objective is to study the relation between the scaled failure rate P_G in (18) and the original failure rate P_F in (2). Towards this goal, we partition the M -dimensional variation space into a large number of identical hyper-rectangles with the same volume and the scaled failure rate P_G in (18) can be approximated as

$$P_G \approx \sum_k I_{\Omega}[\mathbf{x}^{(k)}] \cdot g[\mathbf{x}^{(k)}] \cdot \Delta \mathbf{x}, \quad (20)$$

where $\Delta \mathbf{x}$ denotes the volume of a hyper-rectangle. The approximation in (20) is accurate, if each hyper-rectangle is sufficiently small. Given the definition of $I_{\Omega}(\mathbf{x})$ in (19), Eq. (20) can be re-written as

$$P_G \approx \sum_{k \in \Omega} g[\mathbf{x}^{(k)}] \cdot \Delta \mathbf{x}, \quad (21)$$

where $\{k; k \in \Omega\}$ represents the set of all hyper-rectangles that fall into the failure region.

Substituting (17) into (21), we have

$$P_G \approx \left[\Delta \mathbf{x} / (\sqrt{2\pi} \cdot s)^M \right] \cdot \sum_{k \in \Omega} \exp \left[-\|\mathbf{x}^{(k)}\|_2^2 / 2s^2 \right]. \quad (22)$$

Taking the logarithm on both sides of (22) yields

$$\log P_G \approx \log \frac{\Delta \mathbf{x}}{(2\pi)^{M/2}} - M \cdot \log s + \text{lse}_{k \in \Omega} \left[-\|\mathbf{x}^{(k)}\|_2^2 / 2s^2 \right], \quad (23)$$

where

$$\text{lse}_{k \in \Omega} \left[-\|\mathbf{x}^{(k)}\|_2^2 / 2s^2 \right] = \log \left\{ \sum_{k \in \Omega} \exp \left[-\|\mathbf{x}^{(k)}\|_2^2 / 2s^2 \right] \right\} \quad (24)$$

stands for the log-sum-exp function. The function $\text{lse}(\bullet)$ in (24) is also known as the ‘‘soft maximum’’ from the mathematics [21]. Namely, it is considered as a good approximation of the maximum operator

$$\text{lse}_{k \in \Omega} \left[-\|\mathbf{x}^{(k)}\|_2^2 / 2s^2 \right] \approx \max_{k \in \Omega} \left[-\|\mathbf{x}^{(k)}\|_2^2 / 2s^2 \right]. \quad (25)$$

Substituting (25) into (23) yields

$$\log P_G \approx \alpha + \beta \cdot \log s + \frac{\gamma}{s^2}, \quad (26)$$

where

$$\begin{aligned} \alpha &= \log \frac{\Delta \mathbf{x}}{(2\pi)^{M/2}} \\ \beta &= -M \\ \gamma &= \max_{k \in \Omega} \left[-\|\mathbf{x}^{(k)}\|_2^2 / 2 \right] \end{aligned}. \quad (27)$$

Eq. (26) reveals the important relation between the scaled failure rate P_G and the scaling factor s . The approximation in (26) does not rely on any specific assumption of the failure region. It is valid, even if the failure region is non-convex or discontinuous.

While (27) shows the theoretical definition of the model coefficients α , β and γ , finding their exact values is not trivial. For instance, the coefficient γ is determined by the hyper-rectangle that falls into the failure region Ω and is closest to the origin $\mathbf{x} = \mathbf{0}$. In practice, without knowing the failure region Ω , we cannot directly find out the value of γ . For this reason, we fit the analytical model in (26) by linear

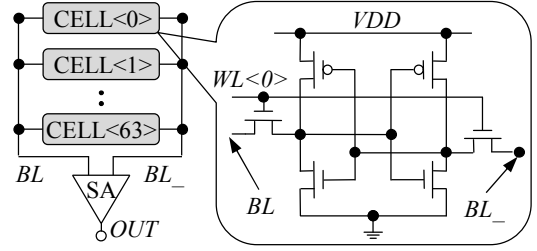


Fig. 5. The simplified schematic is shown for an SRAM column consisting of 64 bit-cells and a sense amplifier (SA) designed in a 45nm CMOS process.

regression. Namely, we first estimate the scaled failure rates $\{P_{G,q}; q = 1, 2, \dots, Q\}$ by setting the scaling factor s to a number of different values $\{s_q; q = 1, 2, \dots, Q\}$. As long as the scaling factors $\{s_q; q = 1, 2, \dots, Q\}$ are sufficiently large, the scaled failure rates $\{P_{G,q}; q = 1, 2, \dots, Q\}$ are large and can be accurately estimated with a small number of random samples. Next, the model coefficients α , β and γ are fitted by linear regression based on the values of $\{(s_q, P_{G,q}); q = 1, 2, \dots, Q\}$. Once α , β and γ are known, the original failure rate P_F in (2) can be predicted by *extrapolation*. Namely, we substitute $s = 1$ into the analytical model in (26)

$$\log P_F^{\text{SSS}} = \alpha + \gamma, \quad (28)$$

where P_F^{SSS} denotes the estimated value of P_F by SSS. Apply the exponential function to both sides of (28) and we have

$$P_F^{\text{SSS}} = \exp(\alpha + \gamma). \quad (29)$$

To make the SSS method of practical utility, maximum likelihood estimation is applied to fit the model coefficients in (26). In addition, a bootstrapping based technique is developed to accurately estimate the confidence interval of the SSS estimator. Due to the page limit, all these implementation details are omitted in this paper, but can be found in the recent publication [18].

Note that to apply SSS, we only need a set of scaling factors and their corresponding scaled failure rates: $\{(s_q, P_{G,q}); q = 1, 2, \dots, Q\}$. As long as $\{s_q; q = 1, 2, \dots, Q\}$ are sufficiently large, $\{P_{G,q}; q = 1, 2, \dots, Q\}$ are not small probability values and, therefore, can be efficiently estimated by CMC. When applying CMC, we only need to determine whether the random samples belong to the failure region. Namely, the PoI does not have to be continuous. Due to this reason, SSS can be applied to estimate the rare failure rates for both continuous and discrete PoIs. However, since SUS explores additional information from the continuous performance values, SUS is often preferred over SSS when we handle continuous PoIs.

To demonstrate the efficacy of SSS, we consider an SRAM column consisting of 64 bit-cells and a sense amplifier (SA) designed in a 45nm CMOS process. Fig. 5 shows the simplified circuit schematic of this SRAM column example. Similar to the SRAM read current example shown in Fig. 2, we consider the local V_{TH} mismatch of each transistor as an independent Normal random variable. In total, we have 384 independent random variables. In this example, the output of SA is considered as the PoI. If the output is correct, we consider the circuit as ‘‘PASS’’. Hence, the PoI is binary, and we cannot apply SUS in this example. For comparison purposes, we run MNIS [10] and SSS for 100 times with 6000 transistor-level simulations in each run. As shown in Fig. 6,

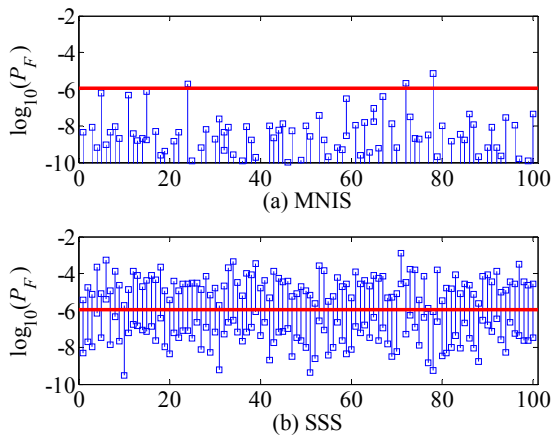


Fig. 6. The 95% confidence intervals (blue bars) of the SRAM example are estimated from 100 repeated runs with 6000 transistor-level simulations in each run for: (a) MNIS and (b) SSS. The red line represents the “golden” failure rate.

there are 3 and 97 CIs out of 100 runs that cover the “golden” failure rate for MNIS and SSS, respectively. Here, the “golden” failure rate is estimated by CMC with 10^9 random samples. MNIS, again, fails to accurately estimate the corresponding CIs. SSS, however, successfully estimates the CIs. These results demonstrate that SSS is superior to the traditional MNIS method in this SRAM example, where the dimensionality of the variation space is more than a few hundred.

IV. Conclusions

Rare failure event analysis in a high-dimensional variation space has attracted more and more attention due to aggressive technology scaling. To address this technical challenge, we summarize two novel approaches: SUS and SSS. Several SRAM examples are used to demonstrate the efficacy of SUS and SSS. More experimental results of SUS and SSS can be found in the recent publications [18]-[19]. Both SUS and SSS are based upon solid mathematical background and do not pose any specific assumption on the failure region. Hence, they can be generally applied to estimate the rare failure rates of a broad range of other circuits, e.g., DFF.

Acknowledgements

This work has been supported in part by the National Science Foundation under contract CCF-1016890 and CCF-1148778.

References

- [1] B. Calhoun, et al., “Digital circuit design challenges and opportunities in the era of nanoscale CMOS,” *Proc. IEEE*, vol. 96, no. 2, pp. 343-365, Feb. 2008.
- [2] A. Bhavnagarwala, X. Tang and J. Meindl, “The impact of intrinsic device fluctuations on CMOS SRAM cell stability,” *IEEE JSSC*, vol. 36, no. 4, pp. 658-665, Apr. 2001.
- [3] R. Heald and P. Wang, “Variability in sub-100nm SRAM designs,” *IEEE ICCAD*, pp. 347-352, 2004.
- [4] R. Kanj, R. Joshi and S. Nassif, “Mixture importance sampling and its application to the analysis of SRAM designs in the presence of rare failure events,” *IEEE DAC*, pp. 69-72, 2006.
- [5] R. Topaloglu, “Early, accurate and fast yield estimation

- through Monte Carlo-alternative probabilistic behavioral analog system simulations,” *IEEE VTS*, pp. 137-142, 2006.
- [6] C. Gu and J. Roychowdhury, “An efficient, fully nonlinear, variability aware non-Monte-Carlo yield estimation procedure with applications to SRAM cells and ring oscillators,” *IEEE ASP-DAC*, pp. 754-761, 2008.
- [7] L. Dolecek, M. Qazi, D. Shah and A. Chandrakasan, “Breaking the simulation barrier: SRAM evaluation through norm minimization,” *IEEE ICCAD*, pp. 322-329, 2008.
- [8] J. Wang, S. Yaldiz, X. Li and L. Pileggi, “SRAM parametric failure analysis,” *IEEE DAC*, pp. 496-501, 2009.
- [9] A. Singhee and R. Rutenbar, “Statistical blockade: very fast statistical simulation and modeling of rare circuit events, and its application to memory design,” *IEEE TCAD*, vol. 28, no. 8, pp. 1176-1189, Aug. 2009.
- [10] M. Qazi, M. Tikekar, L. Dolecek, D. Shah and A. Chandrakasan, “Loop flattening and spherical sampling: highly efficient model reduction techniques for SRAM yield analysis,” *IEEE DATE*, pp. 801-806, 2010.
- [11] R. Fonseca, L. Dilillo, A. Bosio, P. Girard, S. Pravossoudovitch, A. Virazel and N. Badereddine, “A statistical simulation method for reliability analysis of SRAM core-cells,” *IEEE DAC*, pp. 853-856, 2010.
- [12] K. Katayama, S. Hagiwara, H. Tsutsui, H. Ochi and T. Sato, “Sequential importance sampling for low-probability and high-dimensional SRAM yield analysis,” *IEEE ICCAD*, pp. 703-708, 2010.
- [13] S. Sun, Y. Feng, C. Dong and X. Li, “Efficient SRAM failure rate prediction via Gibbs sampling,” *IEEE TCAD*, vol. 31, no. 12, pp. 1831-1844, Dec. 2012.
- [14] R. Kanj, R. Joshi, Z. Li, J. Hayes and S. Nassif, “Yield estimation via multi-cones,” *IEEE DAC*, pp. 1107-1112, 2012.
- [15] S. Au and J. Beck, “Estimation of small failure probabilities in high dimensions by subset simulation,” *Probabilistic Eng. Mechanics*, vol. 16, no. 4, pp. 263-277, Oct. 2001.
- [16] A. Guyader, N. Hengartner and E. Matzner-Löber, “Simulation and estimation of extreme quantiles and extreme probabilities,” *Appl. Math. Optimization*, vol. 64, no. 2, pp. 171-196, Oct. 2011.
- [17] F. Cérou, P. Moral, T. Furon and A. Guyader, “Sequential Monte Carlo for rare event estimation,” *Stat. Computing*, vol. 22, no. 3, pp. 795-808, May 2012.
- [18] S. Sun, X. Li, H. Liu, K. Luo and B. Gu, “Fast statistical analysis of rare circuit failure events via scaled-sigma sampling for high-dimensional variation space,” *IEEE ICCAD*, pp. 478-485, 2013.
- [19] S. Sun and X. Li, “Fast statistical analysis of rare circuit failure events via subset simulation in high-dimensional variation space,” *IEEE ICCAD*, 2014.
- [20] C. Bishop, *Pattern Recognition and Machine Learning*, Prentice Hall, 2007.
- [21] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge University Press, 2009.