

# Accurate and Efficient Gate-Level Parametric Yield Estimation Considering Correlated Variations in Leakage Power and Performance

Ashish Srivastava, Saumil Shah, Kanak Agarwal, Dennis Sylvester, David Blaauw, Stephen Director  
University of Michigan, EECS Department, Ann Arbor, MI 48109  
{ansriv, saumil, agarwalk, dennis, blaauw, director}@eecs.umich.edu

## Abstract

Increasing levels of process variation in current technologies have a major impact on power and performance, and result in parametric yield loss. In this work we develop an efficient gate-level approach to accurately estimate the parametric yield defined by leakage power and delay constraints, by finding the joint probability distribution function (jpdf) for delay and leakage power. We consider inter-die variations as well as intra-die variations with correlated and random components. The correlation between power and performance arise due to their dependence on common process parameters and is shown to have a significant impact on yield in high-frequency bins. We also propose a method to estimate parametric yield given the power/delay jpdf that is much faster than numerical integration with good accuracy. The proposed approach is implemented and compared with Monte Carlo simulations and shows high accuracy, with the yield estimates achieving an average error of 2%.

**Categories and Subject Descriptors:** Performance Analysis and Design Aids

**General Terms:** Reliability, Performance

**Keywords:** Yield, Variability, Leakage, Correlation

## 1. Introduction and Overview of Approach

Process variability has grown in recent technologies due to random dopant effects in small devices, the patterning of features smaller than the wavelength of the optical lithography system and related trends. These variations have a tremendous impact on both power and performance of current integrated circuit (IC) designs. In particular, leakage power has grown to contribute a significant fraction of total power and is known to be highly susceptible to process variations due to its exponential dependence on threshold voltage [1]. In [2], a 20X variation in leakage power for 30% delay variation between fast and slow dies is reported. Due to the inverse correlation of power and delay, most of the fastest chips in a lot are found to have unacceptable leakage and vice versa. This leads to a two-sided constraint on the feasible region in the parametric space and results in significant parametric yield loss.

This yield loss will worsen in future technologies due to increasing process variations and the continued significance of leakage power. Another troublesome observation is that increased variation not only results in a larger spread of leakage power but also in higher average leakage power. Most current optimization approaches do not consider process variations and are unaware of their impact on yield. These approaches invariably result in the formation of a timing wall and result in yield loss due to increased susceptibility to process variations [3].

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

*DAC 2005*, June 13–17, 2005, Anaheim, California, USA.

Copyright 2005 ACM 1-59593-058-2/05/0006...\$5.00.

Several approaches were recently proposed to perform statistical timing or power optimization [4-6], however these approaches neglect the correlation of power and performance, and therefore performing timing yield optimization results in yield loss due to the power constraint while power minimization techniques will harm timing-based yield. Hence, there is a critical need to develop accurate and computationally efficient yield estimation approaches in order to enable yield-driven optimization tools.

Previous work in yield estimation has been limited to predicting either timing [7-12] or (leakage) power yield [13]. Recently [14] presented a chip-level approach to estimate the yield in separate frequency bins given a power constraint. This high-level approach is based on global circuit parameters such as the total device width on a chip. Since it does not use circuit specific information from a gate-level netlist, it is difficult to use for optimization of gate-level parameters, such as the threshold voltage and sizes of individual gates. Another important requirement for an accurate yield estimation approach is to consider all classes of variations which have significantly different impact on delay [15] and power [13]. Process variations are typically classified into inter-die and intra-die components. Intra-die variations are further classified as having correlated and random components. Traditionally inter-die variations have been the dominant source of variations but with process scaling, the random and correlated components of intra-die variations now exceed inter-die variations [16]. The relative magnitude of these components of variation also depends on the process parameter being considered. For example, gate-length variations are generally considered to have roughly comparable random and correlated components whereas gate-length independent threshold voltage is commonly assumed to vary randomly due to random dopant fluctuations [17].

In this work we propose a novel approach to compute the parametric yield of a circuit with high efficiency and accuracy given leakage power and delay constraints. To the best of our knowledge, this is the first such gate-level approach to estimating parametric yield. The variations in process parameters are assumed to be normally distributed. We model the correlated components using a principal component based approach that allows us to express the underlying variations in terms of independent Gaussian random variables (RVs) and consider both inter-die variations and the correlated component of intra-die variations. We perform statistical timing analysis in the spirit of [7] with an additional random component of variation that is not explicitly considered in [7]. Since it is impractical to maintain a separate RV for the random component associated with each gate in a circuit, we lump all random variations into one additional term.

We then develop an approach to perform statistical leakage power analysis and express the circuit leakage power in terms of the same underlying process variations used to express the delay of the circuit. With increasing circuit size the impact of the random component of variation on the variance of power reduces to zero due to the central limit theorem [18]. Thus, although the random component impacts the statistical circuit delay, the correlation between the random components of power and delay has a vanishingly small impact on

**Table 1. Estimated yield for different values of correlation coefficients. Power constraint is set at 1.5X nominal power**

	Estimated Yield		
	Corr = -1.0	Corr = 0.0	Corr = 1.0
Yield	max	$(0.5+\Phi(D))^*$	$0.5+$
Expression	$(\Phi(D)+\Phi(P),0)$	$(0.5+\Phi(P))$	$\Phi(\min(D,P))$
D=-1	0.000	0.095	0.159
D=0	0.100	0.300	0.500
D=1	0.441	0.505	0.600
D=2	0.577	0.586	0.600
D=3	0.599	0.599	0.600

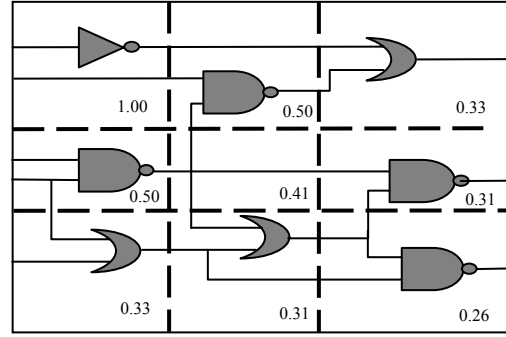
the overall correlation in power and delay. Our results show that even for small circuits with a few hundred gates, the random component has negligible impact on the overall variance of power. This allows us to calculate the correlation in power and delay, and to construct their joint probability distribution function (jpdf). Based on this jpdf of delay and power we develop a closed-form approach to estimate the yield of a design given the delay and leakage constraints. To demonstrate the importance of the power/delay correlation, Table 1 shows yield for varying values of correlation factors at which simple expressions for yield can be obtained.  $\Phi(x)$  represents the value of the Gaussian cumulative distribution function at  $x$ . The yields are estimated for delay constraints of ‘D’ standard deviations (SD) from the mean at a fixed power constraint ‘P’. The results in Table 1 clearly show that the correlation of power and delay has a strong impact on parametric yield, particularly for mid- to high-performance speed bins.

The remainder of the paper is organized as follows. Section 2 briefly reviews the principal component based approach to model process variations. Section 3 presents the core statistical power and timing analysis approaches. In Section 4 we develop an approach to estimate the yield given power and delay constraints. Section 5 presents the results and compares our approach to Monte Carlo circuit simulations. We conclude in Section 6.

## 2. Modeling Process Variations

This section details the variability modeling infrastructure used in our work. Much of this framework is similar in spirit to [7,9] for statistical timing analysis – we also use the same models to consider leakage variability such that the correlation between power and delay is preserved for yield estimation.

In this paper we consider process variations in gate length and gate length-independent threshold voltage ( $V_{th0}$ ) although the approach can be easily extended to consider other sources of variations. The process parameters are expressed as a sum of correlated and random components and the sum of variances of both these components provides the overall variation in the process parameter. To handle the correlated components of variations (inter-die and correlated intra-die) the overall chip area is divided into a grid as shown in Figure 1. In the absence of inter-die process variations the correlation coefficient varies from one (within the same square of the grid) and falls off to zero with increasing distance. Due to inter-die process variations squares on the grid that lie at the opposite corners are correlated and the correlation will fall off to a value higher than zero that depends on the relative contribution of inter-die variations to the correlated variation. Each square in the grid corresponds to a RV of the process parameter which has correlations with all other RVs corresponding to other squares in the grid. The values at the bottom right corner of each of the grids show the correlation coefficients with the top-left square on the grid. Squares that are much further apart demonstrate lower correlation compared to that of adjacent squares on the grid in this model.



**Figure 1. Example partition of a circuit using a grid to model the correlated component of variation (correlation coefficients referenced to top-left element)**

To simplify the problem, this set of correlated RVs is replaced by another set of mutually independent RVs with zero mean and unit variance using the principal components of the set of correlated RVs. A vector of RVs (say  $X$ ) with a correlation matrix  $C$ , can be expressed as a linear combination of the principal components  $Y$  as:

$$X = \Delta_x + \Omega V^{-1} D^{1/2} Y \quad (1)$$

where  $\Delta_x$  is the vector of the mean values of  $X$ ,  $\Omega$  is a diagonal matrix with the diagonal elements being the standard deviations of  $X$ ,  $V$  is the matrix of the eigenvectors of  $C$ , and  $D$  is a diagonal matrix of the eigenvalues of  $C$ . Since the correlation matrix of a multivariate (nondegenerate) Gaussian RV is positive-definite, all elements of  $D$  are positive and the square-root in (1) can be evaluated. We express the delay and leakage power<sup>1</sup> of an individual gate as shown in (2):

$$Delay = d_{nom} + \sum_{i=1}^p \alpha_p (\Delta P_p) \quad (2)$$

$$Leakage = \exp \left( V_{nom} + \sum_{i=1}^p \beta_p (\Delta P_p) \right)$$

where  $d_{nom}$  and  $\exp(V_{nom})$  are the nominal values of delay and leakage power respectively, and the  $\alpha$ 's and  $\beta$ 's represent the sensitivities of delay and the log of leakage to the process parameters under consideration. The variable  $\Delta P_p$  represents the change in the process parameters from their nominal value.

In a statistical scenario, the process parameters are modeled as RVs. If the overall circuit is partitioned using a grid as shown in Figure 1, the delay of individual gates can be expressed as a function of these RVs. Using the principal component approach, the delay in (2) can then be expressed as:

$$Delay = d_{nom} + \sum_{i=1}^p \left( \alpha_p \sum_{j=1}^n \gamma_{ji} z_j \right) + \eta_d R \quad (3a)$$

where  $z_j$ 's are the principal components of the correlated RV's  $\Delta P_p$ 's in (2) and the  $\gamma$ 's can be obtained from (1).  $R \sim N(0,1)$  in the above equation represents the random component of the variations of all process parameters lumped into a single term that contributes a total variance of  $\eta_d^2$  to the overall variance of delay. Similarly the leakage power for an individual gate can be expressed as:

$$Leakage = \exp \left( V_{nom} + \sum_{i=1}^p \left( \beta_p \sum_{j=1}^n \gamma_{ji} z_j \right) + \eta_l R \right) \quad (3b)$$

<sup>1</sup> Since variations in dynamic power are typically much smaller than those observed in static power, we focus on statistical leakage power analysis for yield estimation in this paper.

The next section shows that these representations of delay and power allow for significant simplification in the joint timing and power analysis, which otherwise becomes computationally inefficient if the spatial correlation is maintained without such a unified principal component approach.

### 3. Statistical Analysis

In this section we first provide an overview of the statistical timing analysis in [7] which we have extended to consider both correlated and random components of variations. We then provide the details of the approach for performing statistical leakage power analysis. During timing analysis the arrival time at each node is maintained in the same canonical form as the delay of the individual gates, which enables an efficient approach for the traversal of the timing graph. Similarly, during power analysis the sum of leakage power is maintained in a canonical form as the leakage of different gates is summed.

#### 3.1 Timing Analysis

The delay of each gate (say  $a$ ) can be expressed as follows using the expression developed in the previous section:

$$D_a = a_0 + \sum_{i=1}^n a_i z_i + a_{n+1} R \quad (4)$$

This serves as the canonical expression for delay. The mean delay is simply the nominal delay ( $a_0$ ). Since the principal components ( $z_i$ 's) are uncorrelated  $N(0,1)$  RVs, the variance of the delay can be expressed as:

$$Var(D_a) = \sum_{i=1}^n (a_i^2) + a_{n+1}^2 \quad (5)$$

and the covariance of the delay with one of the principal components can be obtained as:

$$Cov(D_a, z_i) = E(D_a z_i) - E(D_a)E(z_i) = a_i^2 \quad \forall i \in \{1, 2, \dots, n\} \quad (6)$$

In [9], it was shown that delay distributions arising due to correlated reconvergent fanouts can be tightly upper-bounded by assuming them to be independent. Since the random components are uncorrelated and do not contribute to the covariance of the delay at the two nodes at the input of a gate (e.g., ' $a$ ' and ' $b$ '), the covariance can be obtained as:

$$Cov(D_a, D_b) = \sum_{i=1}^n a_i b_i \quad (7)$$

In deterministic timing analysis the delay of the circuit is found by applying two functions to the delay of individual gates: sum and max. Similar functions for the canonical delay expressions are defined as:

$$Sum(D_a, D_b) = (a_0 + b_0) + \sum_{i=1}^n (a_i + b_i) z_i + \sqrt{a_{n+1}^2 + b_{n+1}^2} R \quad (8)$$

The max function of normally distributed RVs is not a strict Gaussian. References [7,19] have shown that the maximum of two Gaussian RVs can be closely approximated by another Gaussian. If ' $c$ ' is the max of ' $a$ ' and ' $b$ ' where  $a$  and  $b$  are Gaussian RV's, then the parameters of  $c$ , which is assumed to be Gaussian, can be obtained using expressions developed in [20]. This approach provides the mean and variance of  $c$  in terms of the mean and variance of  $a$  and  $b$  and their correlation coefficient. Reference [20] also develops expressions to evaluate the correlation of  $c$  with any other RV in terms of the correlation of the RV with  $a$  and  $b$ . In the spirit of [7,8] we assume that  $c$  can again be expressed in the same canonical form as  $a$  and  $b$ . To find the coefficients in the expression for  $c$  in canonical form, the mean, variance and the correlations of  $c$  with the principal components are matched, giving

$$c_0 = E(\max(a, b)) \quad (9)$$

$$c_i = \text{cov}(c, z_i) = \text{cov}(\max(a, b), z_i) \quad \forall i \in \{1, \dots, n\}$$

$$c_{n+1} = \left( Var(\max(a, b)) - \sum_{i=1}^n c_i^2 \right)^{1/2}$$

By modeling the random component, we can preserve the mean, variance, and correlations, avoiding the need to scale the coefficients of the principal components to match variance, which loses their correlation [7]. To compute the max of more than two variables the above technique is applied recursively.

Using the timing analysis approach outlined above, we can develop an expression for the delay of a circuit in terms of the RVs associated with process parameter variations. In the next sub-section we develop an approach for gate-level statistical leakage power analysis. The key in this step is to preserve the correlation in delay and power, which is achieved by using a similar principal component based approach with the same underlying RVs.

#### 3.2 Power Analysis

We express the leakage power of each gate as a lognormal (exponential of a Gaussian) RV based on the power model discussed in Section 2. The leakage power of the total circuit can then be expressed as a sum of correlated RVs. This sum can be accurately approximated as another lognormal random variable [21]. Reference [21] shows that the approximation performed using an extension of Wilkinson's method [22] (based on matching the first two moments) provides good accuracy. The leakage power of an individual gate  $a$  is expressed as

$$P_{leak}^a = \exp\left(a_0 + \sum_{i=1}^n a_i z_i + a_{n+1} R\right) \quad (10)$$

where the  $z$ 's are principal components of the RVs and the  $a$ 's are the coefficients obtained using (1) and (2). The mean and variance of the RV in (10) can then be computed as

$$E(P_{leak}^a) = \exp\left(a_0 + \frac{1}{2} \sum_{i=1}^{n+1} a_i^2\right) \quad (11)$$

$$Var(P_{leak}^a) = \exp\left(2a_0 + \sum_{i=1}^{n+1} a_i^2\right) - \exp\left(2a_0 + \frac{1}{2} \sum_{i=1}^{n+1} a_i^2\right) \quad (12)$$

The correlation of the leakage of gate  $a$  with the lognormal RV associated with  $z_j$  is found by evaluating

$$E\left(P_{leak}^a e^{z_j}\right) = \exp\left(a_0 + \frac{1}{2} \sum_{i=1, i \neq j}^{n+1} a_i^2 + (a_j + 1)^2\right) \quad \forall j \in \{1, 2, \dots, n\} \quad (13)$$

Similarly the covariance of the leakage of two gates ( $a$  and  $b$ ) can be obtained by using

$$E\left(P_{leak}^a P_{leak}^b\right) = \exp\left((a_0 + b_0) + \frac{1}{2} \left(\sum_{i=1}^n (a_i + b_i)^2 + a_{n+1}^2 + b_{n+1}^2\right)\right) \quad (14)$$

We assume that the sum of leakage power can be expressed in the same canonical form as (10). If the random variables associated with all the gates in the circuit are summed in a single step the overall complexity of the approach is  $O(n^2)$  due to the size of the correlation matrix. Since the sum of two lognormal RVs is assumed to have a lognormal distribution in the same canonical form, we can use a recursive technique to estimate the sum of more than two lognormal RVs. In each recursive step we sum two RVs of the form in (10) to obtain another RV in the same canonical form. To find the coefficients in the expression for the sum of the RVs we match the first two moments (as in Wilkinson's method) and the correlations with the lognormal RVs associated with each of the Gaussian principal components. We now outline one of the recursive steps where we sum  $P_{leak}^b$  and  $P_{leak}^c$  to obtain  $P_{leak}^a$ .

The coefficients associated with the principal components can be found using (11-14) and expressing the coefficients associated with the principal components as

$$a_i = \log \left( \frac{E(P_{leak}^a e^{z_i})}{E(P_{leak}^a)E(e^{z_i})} \right) = \log \left( \frac{E(P_{leak}^b e^{z_i}) + E(P_{leak}^c e^{z_i})}{(E(P_{leak}^b) + E(P_{leak}^c))E(e^{z_i})} \right) \quad (15)$$

Using the expressions developed in [13], the remaining two coefficients in the expression for  $P_{leak}^a$  can be expressed as

$$a_0 = \frac{1}{2} \log \left( \frac{(E(P_{leak}^b) + E(P_{leak}^c))^4}{(E(P_{leak}^b) + E(P_{leak}^c))^2 + Var(P_b) + Var(P_c) + 2Cov(P_b, P_c)} \right) \quad (16)$$

$$a_{n+1} = \left[ \log \left( 1 + \frac{Var(P_b) + Var(P_c) + 2Cov(P_b, P_c)}{(E(P_{leak}^b) + E(P_{leak}^c))^2} \right) - \sum_{i=1}^n a_i^2 \right]^{0.5} \quad (17)$$

Having obtained the sum of two lognormals in the original canonical form, the process is recursively repeated to compute the expression for the total leakage power of the circuit.

The timing and power analysis techniques outlined in this section is used to efficiently estimate the individual probability distribution functions of delay and power. The correlation in delay and leakage power arising from the correlated components of variation can be easily estimated since the correlated variations are expressed in terms of the principal components used to develop the expressions for both delay and power. As will be shown in the results section, the dependence of the variance of leakage power on the random component is very weak. This arises due to the fact that the random component associated with each gate is independent and hence the ratio of standard deviation to mean for the sum of these independent RVs is inversely proportional to the square root of the number of RVs summed [18]. This ratio does not reduce for correlated RVs – therefore if a large number of RVs are summed with both correlated and random components, the overall variance is dominated by the variance of the correlated component. Hence, the correlation due to the random component, which is difficult to compute efficiently, is also insignificant and can be safely neglected.

#### 4. Yield Estimation

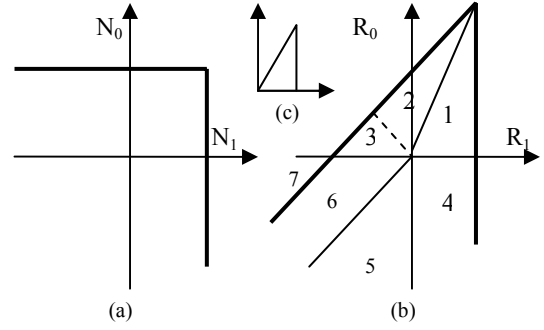
The parametric yield of a circuit given delay and power constraints can be expressed as

$$Y = P(D \leq D_0, P \leq P_0) \quad (18)$$

which is the probability of the circuit delay and power being less than  $D_0$  and  $P_0$  respectively. Since delay and power are correlated, the yield cannot be simply computed by multiplying the separate probabilities. To express the yield as the probability of a bivariate Gaussian RV we take the logarithm of the leakage power constraint. The correlation coefficient of the two Gaussian RVs in the yield equation can now be obtained using (7). We express the yield in terms of two  $N(\theta, I)$  RVs  $N_0$  and  $N_1$  as

$$Y = P \left( N_0 \leq \frac{D_0 - \mu_D}{\sigma_D}, N_1 \leq \frac{\log(P_0) - \mu_{\log(P)}}{\sigma_{\log(P)}} \right) \quad (19)$$

Since the correlation coefficient of two RVs does not change under a linear transformation with positive coefficients of the original RVs, the correlation between  $N_0$  and  $N_1$  remains same. One approach to evaluate this expression is to perform numerical integration of the jpdf over the feasible region, but this makes the approach computationally inefficient. A look-up table based approach, though efficient, involves substantial inaccuracy due to the required interpolation as noted in [23]. Hence, we adopt an analytical approach to estimate the yield which makes the approach efficient and practical within a yield optimization framework. The feasible region defined using two correlated RVs is transformed to a set of two uncorrelated RVs using the following transformation:



**Figure 2. Transformation of the feasible region from (a) to (b) under the transformation expressed in (20)**

$$R_0 = N_0; R_1 = \left( \frac{N_1 - \rho N_0}{(1 - \rho^2)^{1/2}} \right) \quad (20)$$

This transformation maps the feasible region from a rectangle to a triangle as shown in Figure 2 for the case where  $\rho < 0$ , which is the case of interest. The desired probability can be obtained by using approximate expressions developed in [23] for evaluating probabilities of uncorrelated bivariate Gaussian RVs in regions of the form shown in Figure 2(c). To evaluate the probability of the region shown in Figure 2(b) we partition the figure as shown. The desired probability can then be expressed as a sum of the probabilities in Regions 1-6 which can be evaluated as follows:-

*Region 1:* Already in the form required in [23],

*Region 2:* Since the integral of the region is circularly symmetric, if the axes are rotated such that the dotted line as shown in Figure 2(b) lies along the x-axis then Region 2 is again in the same form as Figure 2(c),

*Region 4:* The probability in this region is

$$P(R_0 \leq 0, 0 \leq R_1 \leq X) \quad (21)$$

where  $X$  is the point where the vertical line cuts the  $R_0$  axis. Since  $R_0$  and  $R_1$  are statistically independent, this probability can be simply expressed as

$$P(R_0 \leq 0)P(0 \leq R_1 \leq X) = 0.5\phi(X) \quad (22)$$

where  $\Phi$  is the normal integral from  $0$  to  $x$ ,

*Regions 3, 5 and 6:* The probability for these regions can be expressed as

$$P(R_0 \leq 0, R_1 \leq 0) + P(3+6) - P(6+7) \quad (23)$$

The first and second terms in (23) correspond to a region that has the same form as Region 4 and the region for the third term has the same form as Region 1. Thus, the desired yield expressed in (18) can be efficiently estimated using closed-form expressions.

In terms of computational complexity, the proposed approach differs from [7] in the computation of an extra term associated with the random component. Thus the overall complexity of the timing analysis remains  $O(nN_g)$ , where  $n$  is the number of terms in the delay expression that corresponds to the number of partitions into which the circuit is divided, and  $N_g$  is the number of gates in the circuit. The power analysis is similar and requires an additional  $O(nN_g)$  steps. The correlation computation requires an additional  $O(n)$  steps, and the yield estimation runs in constant time. The computation of the principal components requires  $O(pn^3)$  steps where  $p$  is the number of process parameters required. The cubic dependence results from the eigenvector computation required



