

Joint Design-Time and Post-Silicon Optimization for Digitally Tuned Analog Circuits

Wei Yao, Yiyu Shi, Lei He and Sudhakar Pamarti

Electrical Engineering Dept., University of California, Los Angeles, CA, 90095, USA
 {weiyao, yshi, lhe, spamarti}@ee.ucla.edu

Abstract— Joint design time and post-silicon optimization for analog circuits has been an open problem in literature because of the complex nature of analog circuit modeling and optimization. In this paper we formulate the co-optimization problem for digitally tuned analog circuits to optimize the parametric yield, subject to power and area constraints. A general optimization framework combining the branch-and-bound algorithm and gradient ascent method is proposed. We demonstrate our framework with two examples in high-speed serial link, the transmitter design and the phase-locked-loop (PLL) design. Simulation results show that compared with the design heuristic from analog designers’ perspective, joint design-time and post-silicon optimization can improve the yield by up to 47% for transmitter design and up to 56% for PLL design under the same area and power constraints. To the best of the authors’ knowledge, this is the first in-depth study on yield-driven analog circuit design technique that optimizes post-silicon tuning together with the design-time optimization.

I. INTRODUCTION

As process technologies scale down to 90nm and below, traditional circuit design methodologies are confronted by the prominent problem of process variation. Conventional corner-based design methodology for analog circuits, in particular, is becoming insufficient and may eventually be inviable as variation increases. Statistical design is therefore proposed and it defines *parametric yield* as the probability that the design meets a specified performance or power constraint. Different techniques exist to maximize the parametric yield for analog circuits and generally fall into two complementary categories: *design-time optimization* and *post-silicon tuning*.

Design-time optimization techniques explore the design space at system-level and device-level to maximize the yield for analog circuits. At system-level, different circuit architectures are explored for a trade-off between power, area, and performance. Architectures such as closed-loop negative feedback have good immunity from process variation. In addition, impacts of process variation can also be reduced by transistor sizing and layout optimization.

Post-silicon tuning in analog design has been widely adopted to combat process variation. Tunable elements such as programmable capacitance array (PCA) and resistance array are proposed to adjust analog circuit performance after chip fabrication [1], [2], [3]. Fig. 1 shows two examples of the tunable elements in analog design: tunable CMOS current source and capacitance array, where β is the resolution (number of control bits). By applying appropriate, potentially different, control signals $D[i]$ ($1 \leq i \leq \beta - 1$) on individual chips, performance can be adjusted to maximize yield. While this will be discussed in more detail in Section II, in both examples the tuning values are digitized. Such *digitally tuned analog*

This paper is partially supported by SRC and a UC MICRO grant sponsored by Actel and Fujitsu. Address comments to lhe@ee.ucla.edu.

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, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ICCADCAD '09, November 2-5, 2009, San Jose, California, USA.
 Copyright 2009 ACM 978-1-60558-800-1/09/11 ...\$10.00.

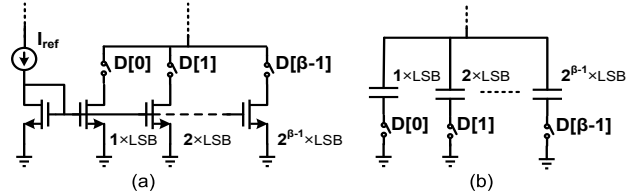


Fig. 1. Examples of digitally tuned analog circuits: (a) CMOS current source and (b) capacitance array.

circuits have wide applications because of their noise-insensitivity and good technology scalability [4].

Post-silicon tuning has been shown to directly impact the design-time optimization for analog circuits [4]. On one hand, post-silicon tunability can significantly relax the analog design by providing a certain capacity to “correct” performance deviation after fabrication. On the other hand, tuning circuitry consumes extra area and power which needs to be considered during design-time optimization in order to meet design specifications. It is natural to expect that by extending joint design-time and post-silicon optimization to analog design, a better parametric yield can be achieved. The complication of modeling and optimizing tunable analog circuits, however, leaves co-optimization an open problem in literature.

In this paper, we study the joint design-time and post-silicon optimization with focus on digitally tuned analog circuits where design variables such as transistor sizes are continuous, while variables such as the tuning resolution are discrete in nature. Assuming that the resolutions are the only changing variables and all the remaining variables are fixed, we derive the performance upper bound among all permissible resolutions. We propose a general optimization framework combining the branch-and-bound algorithm on the resolutions and gradient-ascent method on the unpruned branches. We use the high-speed serial link as our application and provide two analog design examples to demonstrate the joint optimization framework: transmitter equalization filter design and phase-lock loop (PLL) design. Simulation results show that compared with the design heuristic from analog designers’ perspective, joint design-time and post-silicon optimization can improve the yield by up to 47% for transmitter equalization filter design and up to 56% for PLL design under the same area and power constraints.

The remaining of the paper is organized as follows: Section II briefly reviews the post-silicon tuning technique and Section III provides the formulation for our joint optimization problem. Section IV discusses the proposed optimization framework and the designs for transmitter and PLL circuits in high-speed serial link are discussed in Sections IV and V. Experimental results are presented in Section VI and concluding remarks are given in Section VII.

II. PRELIMINARIES ON DIGITALLY TUNED ANALOG CIRCUITS

Analog circuits are very sensitive to process, voltage, and temperature (PVT) variations. Among all sources of variations, the random mismatches caused by doping fluctuations are expected to become dominant within the next few technology generations. In this paper, we focus on the transistor drain current variation caused by threshold voltage (V_{th}) mismatch.

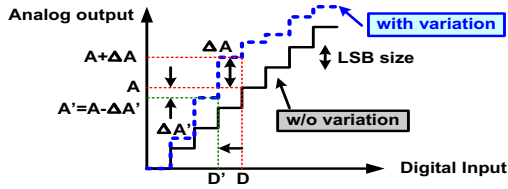


Fig. 2. Post-silicon tuning through DAC

Such drain current variation leads to significant power and performance degradation and various analog design techniques are proposed to reduce its impact. In particular, post-silicon tuning is widely used to calibrate process variation after fabrication using tunable elements. Examples of tunable elements can be found in Fig. 1. Those circuits conceptually operate as digital-to-analog conversion (DAC) circuits. By given a control signal, e.g. D , an analog output, e.g. A , is produced proportionally. There are two major design aspects for digitally tuned analog circuits: least-significant-bit (LSB) size and resolution. The LSB size determines the minimum step in the digital-to-analog conversion. In the CMOS current source shown in Fig. 1(a), for example, it physically represents the drain current for the LSB transistor (I_{LSB}). In the capacitance array shown in Fig. 1(b), it represents the minimum size capacitance (C_{LSB}) in the array. Resolution, on the other hand, is the number of bits used as input control signal. Given the LSB size and resolution, the tuning range can be directly determined. In this paper, we denote the resolution as β and the LSB size as γ .

An example of a digital-to-analog conversion curve is shown in Fig. 2. Assume that digital input D is designed to generate analog output A . With the V_{th} variation, however, the conversion curve becomes nonlinear, and input D generates output with a ΔA deviation with respect to A . To make the analog output closer to the desired value, one can change the input from D to D' and, therefore, a smaller deviation $\Delta A'$ can be obtained. In general, post-silicon tuning is performed by increasing or decreasing the input stepwise to find the minimum deviation.

By applying the tuning technique, the effect of process variation can be significantly reduced. Extra circuits, however, are needed to provide tunability. We assume $D = [100]$ and generate $A = 4 \cdot I_{LSB}$ in Fig. 2. In addition to the required 4 LSB current sources, we need to implement a total of 7 LSB current sources to achieve 3-bit tunability. The resulting extra area increases power and such extra area may be used to design-time optimization such as device sizing or alternative architecture to reduce variation. Therefore, an optimal balance between the performance and area/power cost considering design time and post-silicon tuning must be found.

III. PROBLEM FORMULATION

Without loss of generality, analog design-time optimization can be described to determine the optimal design parameters to maximize the parametric yield, subject to the power and area constraints. Mathematically,

$$(\mathbf{P0}) \quad \max \quad \text{Prob}(F(\mathbf{x}) \leq \bar{f}) \quad (1)$$

$$s.t. \quad \text{Prob}(P(\mathbf{x}) \geq \bar{p}) \leq \epsilon, \quad (2)$$

$$A(\mathbf{x}) \leq \bar{a} \quad (3)$$

$$\mathbf{x}_l \preceq \mathbf{x} \preceq \mathbf{x}_u, \quad \mathbf{x} \in R^k \quad (4)$$

where $F(\cdot)$, $P(\cdot)$, and $A(\cdot)$ represent the functions of performance metric, area, and power, respectively. \bar{f} , \bar{p} and \bar{a} are the upper bounds of the performance metric; power and area given by the design specifications; \mathbf{x} is the vector of length k formed by the design variables with lower bound \mathbf{x}_l and upper bound \mathbf{x}_u given by the design specifications; k is the total number of design variables; ϵ is a small positive number indicating the tolerance for power variation over the upper bound \bar{p} .

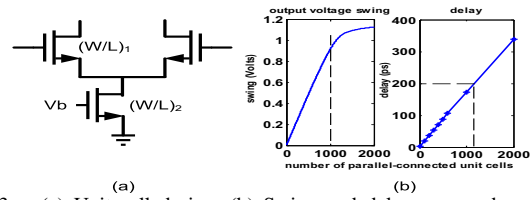


Fig. 3. (a) Unit cell design. (b) Swing and delay vs. number of parallel-connected cells.

With post-silicon tuning, we first consider how to design the digitally tuned elements. In this paper, we adopt a simple but direct method based on the unit cell. An example of a unit cell for the CMOS current source is shown in Fig. 3(a). Assume that we have characterized a total number of m unit cells with different transistor width/length and bias voltage under the condition that they all draw the same amount of current I_{unit} . Each unit cell α_i represents a set of transistor W/L and bias voltage V_b , where $0 \leq \alpha_i \leq m$. Any larger transistor, which draws larger current and provides larger swing at the output, can be obtained by connecting the unit cells of the same type in parallel. Such parallel connection ensures linear relationship for the parasitic capacitance and current driving capability, which is measured by output swing and delay as shown in Fig. 3(b). Moreover, by limiting the maximum number of connected cells, the transistor-level biasing constraints can be guaranteed to ensure all transistors work in the desired operation region. Similar unit cell design methodology can be extended to other digitally tuned elements such as capacitance array.

As a result, the parametric yield can be rewritten as

$$\text{Prob}(\hat{F}(\mathbf{x}, \alpha, \beta, \gamma) \leq \bar{f}) \quad (5)$$

where $\hat{F}(\cdot)$ is the performance metric after tuning, α are the indices of the types of unit cell design, β are vectors representing the resolution used for the digitally tuned elements. γ are the LSB sizes in terms of the number of unit cells used to implement the LSB of the digitally tuned element. The joint design-time and post-silicon optimization then can be extended from $(\mathbf{P0})$ as

$$(\mathbf{P1}) \quad \max \quad \text{Prob}(\hat{F}(\mathbf{x}, \alpha, \beta, \gamma) \leq \bar{f}) \quad (6)$$

$$s.t. \quad \text{Prob}(\hat{P}(\mathbf{x}, \alpha, \beta, \gamma) \geq \bar{p}) \leq \epsilon \quad (7)$$

$$A(\mathbf{x}, \alpha, \beta, \gamma) \leq \bar{a} \quad (8)$$

$$\mathbf{x}_l \preceq \mathbf{x} \preceq \mathbf{x}_u, \quad \mathbf{x} \in R^k \quad (9)$$

$$0 \leq \alpha \leq m\mathbf{1}, \quad \alpha \in Z^n \quad (10)$$

$$0 \leq \beta, \quad \beta \in Z^n \quad (11)$$

$$0 \leq \gamma, \quad \gamma \in Z^n, \quad (12)$$

where m is the total number of unit cell designs, n is the total number of tuning elements in the circuit and $\hat{P}(\cdot)$ is the power consumption after tuning. Note that there is no explicit bound necessary for β and γ as they are implicitly bounded by the power and area constraints (7) and (8).

IV. OPTIMIZATION FRAMEWORK

$(\mathbf{P1})$ is hard to solve because it is a mixed integer non-convex programming problem, the complexity of which grows exponentially with the number of integer variables (the dimension of the vectors α , β and γ). Therefore, we propose to separate the integer variables and the continuous variables. We define a new function $Z(\mathbf{t})$ as the optimum value of $(\mathbf{P1})$ when $\mathbf{x} = \mathbf{t}$. If $(\mathbf{P1})$ is infeasible at $\mathbf{x} = \mathbf{t}$, then $Z(\mathbf{t}) = -\infty$. Accordingly, $\mathbf{P1}$ is equivalent to an unconstrained nonlinear optimization problem with a continuous feasible region:

$$\max Z(\mathbf{t}), \quad \mathbf{t} \in R^k, \quad (13)$$

which can be solved efficiently by the first order gradient method if we can evaluate $Z(\mathbf{t})$ and $\frac{\partial Z(\mathbf{t})}{\partial \mathbf{t}}$ at any point $\mathbf{t} = \hat{\mathbf{t}}$ to find local

maximum. Below we will discuss how to evaluate the function value and first order derivative efficiently.

A. Algorithm Overview

To evaluate $Z(t)$ we need to solve problem **(P1)** for given $x = t$,

$$\text{(P2)} \quad Z(t) = \max \quad \text{Prob}(\hat{F}(t, \alpha, \beta, \gamma) \leq \bar{f}) \quad (14)$$

$$\text{s.t.} \quad \text{Prob}(\hat{P}(t, \alpha, \beta, \gamma) \geq \bar{p}) \leq \epsilon \quad (15)$$

$$A(t, \alpha, \beta, \gamma) \leq \bar{a} \quad (16)$$

$$0 \leq \alpha \leq m\mathbf{1}, \quad \alpha \in \mathbb{Z}^n \quad (17)$$

$$\mathbf{0} \leq \beta, \quad \beta \in \mathbb{Z}^n \quad (18)$$

$$\mathbf{0} \leq \gamma, \quad \gamma \in \mathbb{Z}^n, \quad (19)$$

with variables α , β and γ . **(P2)** is an integer programming problem, which is an NP-hard problem. Though software does exist in literature to solve general integer programming problems, in this paper we propose an optimization framework to efficiently solve it using the special properties of digitally tuned analog circuits.

The optimization framework combines the branch-and-bound (BnB) algorithm with the gradient ascent method (GDA). Assume that we know how to partition the feasibility space into different regions and how to efficiently obtain an upper bound of the objective function (14) for each region. Then, according to the principles of the BnB algorithm, we can prune regions that have an upper bound worse than the existing solutions, thereby maximizing the performance metric. Note that such a pruning technique is exact and not heuristic. If a region cannot be pruned, we employ GDA optimization to find a local maximum in it. The final solution $Z(t)$ is obtained by comparing the optimal solutions found in each unpruned region.

To evaluate the first order derivative $\frac{\partial Z(t)}{\partial t}$, a direct method would be to use the finite different method: Compute $Z(t + \delta e_i)$ ($1 \leq i \leq k$) for some small positive number δ , where e_i is a unit vector with the i^{th} element equal to 1 and other elements equal to 0. Then the i^{th} element of $\frac{\partial Z(t)}{\partial t}$ can be obtained by

$$\frac{\partial Z(t)}{\partial t_i} \approx \frac{1}{\delta} (Z(t + \delta e_i) - Z(t)). \quad (20)$$

As such, the cost for evaluating $\frac{\partial Z(t)}{\partial t}$ would be quite expensive as we would have to solve k integer programming problems. Note that k is the total number of design variables, which can be quite large in practical problems. This urges us to turn to some alternative approach to approximate the computation with affordable cost.

Since we can obtain the upper bound of the objective function in each region efficiently, the upper bound of $Z(t)$ is just the maximum of all those upper bounds. Denoting the upper bound of $Z(t)$ as $\bar{Z}(t)$, the derivative of $Z(t)$ can be approximated by applying finite difference method on $\bar{Z}(t)$, i.e.,

$$\frac{\partial Z(t)}{\partial t_i} \approx \frac{1}{\delta} (\bar{Z}(t + \delta e_i) - \bar{Z}(t)). \quad (21)$$

Note that the accuracy of the approximation depends on how the upper bound is calculated. If the upper bound is tight, then the approximation will converge to the exact derivatives.

B. Partitioning and Bound Estimation

In this paper, we partition the solution space according to the unit cell index and LSB size of each tap. In general, the yield upper bound for a given region is hard to compute. Fortunately, in this particular type of problem, where digitally tuned analog circuits are involved, we are able to obtain the bound through a special relaxation. Suppose we can solve **(P2)** without power and area constraints, then such an optimal value can serve as the upper bound of the constrained problem **(P2)** since we have expanded the feasible space. The physical meaning of solving unconstrained **(P2)** is that given the unit cell design and

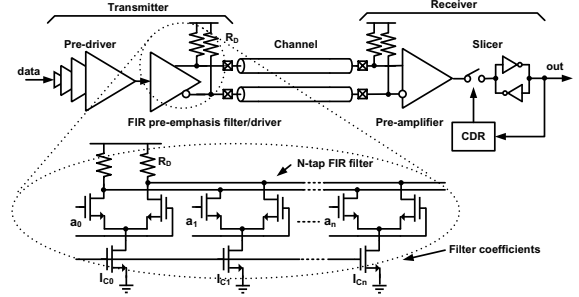


Fig. 4. System diagram of a high-speed serial link.

LSB sizes, find the optimal resolution that gives the maximum yield, which can be solved easily through design-time optimization. The optimality of the solution is guaranteed because any increase in the resolution only increases the total area and the parasitics while the minimum distance to the target values remains the same, which will downgrade the performance.

C. Gradient Ascent Method

Given the partitioning method discussed in the previous section, if a particular region cannot be pruned by comparing its upper bound with the current solution, we need to solve **(P2)** for optimal β with given unit cell indices $\bar{\alpha}$ and LSB size set $\bar{\gamma}$.

In essence, the gradient ascent method sequentially takes steps in a direction proportional to the gradient, until a local maximum of the objective function is reached [5]. At each step we increase/decrease each variable by 1 and check the change of the objective function. Note that by doing so we are actually computing the gradient because all the variables are integers. We then move along the direction that causes the maximum increase. This is iteratively done until the relative change of the objective value is below a certain threshold. The termination of the algorithm indicates that one of the local maxima has been reached or that we have reached the boundary. The initial guess for the GDA can be arbitrarily chosen. In our experiments, we found that it did not influence runtime or quality significantly for both of the examples studied. In addition, we observed that the algorithm always converges to local optimum within two or three iterations.

V. TRANSMITTER DESIGN IN HIGH-SPEED SERIAL LINK

The system diagram of a high-speed serial link is shown in Fig. 4. At the transmitter end, the pre-driver drives the FIR pre-emphasis filter at the designated data rate. The pre-emphasis filter is used to counteract the inter-symbol interference (ISI)[6] caused by the bandwidth-limited channel, which behaves as a transmission line and can be characterized by the Telegrapher's equations with RLGC per-unit-length model. The pre-emphasis filter can be expressed as

$$b_i = \sum_{j=0}^{n-1} W_j a_{i-j}, \quad (22)$$

where n is the number of filter taps, W_i is the tap coefficient for tap i , and a_i is the transmitted non-return-to-zero (NRZ) symbol. At receiver end, the pre-amplifier, along with the slicer decision circuit, is responsible for detecting the data from the received signal. Moreover, the clock is embedded in the transmitted data and the clock data recovery (CDR) sub-system is used to extract the clock from the serial data stream.

In order to focus on the transmitter optimization, in our first example, we assume that the frequency domain response for the channel and the receiver is given. In addition, we assume that an ideal sampling clock is obtained through the receiver CDR circuits.

The performance of the overall system is usually quantified in terms of BER, the rate at which errors occur during data transmission. To start with, we formulate the design-time optimization problem

to minimize the BER of the system subject to power and area constraints. The design variables include the number of taps n of the filter, the transistors sizing W/L , and the bias voltage V_b in the CMOS current source. Assume that we have characterized a total number of m unit cells and each unit cell α_i represents a set of transistor W/L and bias voltage V_b , as shown in Fig. 3.

Since directly measuring the BER requires a long period of time, error vector magnitude (EVM) is used in this paper to estimate the BER because of their monotonic property [7]. We can easily map the EVM to the BER from table look-up and accordingly, the objective function (1) takes the form

$$\text{Prob}(\text{BER}(n, \alpha) \leq \bar{f}). \quad (23)$$

The area $A(n, \alpha)$ and power $P(n, \alpha)$ of the transmitter are mainly contributed by the pre-emphasis filter and the pre-driver, i.e.

$$A(n, \alpha) = A_{\text{pre-driver}}(n, \alpha) + A_{\text{filter}}(n, \alpha), \quad (24)$$

$$P(n, \alpha) = P_{\text{pre-driver}}(n, \alpha) + P_{\text{filter}}(n, \alpha). \quad (25)$$

For tap i ($1 \leq i \leq n$), we use unit cells of type α_i ($1 \leq \alpha_i \leq m$) with the parasitic capacitance $C_{\text{unit}}^{\alpha_i}$ and the occupied area $A_{\text{unit}}^{\alpha_i}$. The required number of cells q_i for that tap is determined by its coefficient W_i and the unit current I_{unit} :

$$q_i = \lceil \frac{W_i}{I_{\text{unit}}} \rceil. \quad (26)$$

Accordingly, the total area used in the pre-emphasis filter can be calculated as

$$A_{\text{filter}}(n, \alpha) = \sum_{i=1}^n q_i A_{\text{unit}}^{\alpha_i}. \quad (27)$$

The total parasitic capacitance C_{para} can be calculated as

$$C_{\text{para}}(n, \alpha) = \sum_{i=1}^n q_i C_{\text{unit}}^{\alpha_i}. \quad (28)$$

The power consumed by the filter (P_{filter}) contains both static power and dynamic switching power and can be expressed as

$$P_{\text{filter}}(n, \alpha) = \rho \sum_{i=1}^n q_i \cdot I_{\text{unit}} \cdot V_{dd} + (1 - \rho) f \cdot V_{dd}^2 \cdot C_{\text{para}}, \quad (29)$$

where f is the data rate. ρ is the ratio between static power and total power, which depends on detailed delay and switching probability and can be obtained from simulation.

The pre-driver is designed according to the total gate capacitance at the filter input $C_{\text{gate}} = \sum_{i=1}^n q_i C_g^{\alpha_i}$, where $C_g^{\alpha_i}$ is the input transistor gate capacitance of unit cell α_i . Here we assume the pre-driver is designed through logic effort using a simple inverter chain and the occupied area and consumed power can be calculated accordingly. Note that other configurations like CML pre-drivers with swing control can also be applied. The optimization problem can then be mathematically formulated as shown in (P0).

In the presence of process variation, assuming transistor threshold voltage V_{th} has a normal distribution with 10% variation [8], the power consumed by the transmitter varies by 30% variation and the BER varies in the magnitude of $10^8 \times$ for the same design, as demonstrated in Fig. 5(a). By applying the tuning technique, simulation results show that the span of power and BER variation becomes much smaller as shown in Fig. 5(b). Extra circuits, however, are needed to provide this tunability and an optimal balance between the performance and area/power cost has to be found.

To cast the problem into the format of (P1), we need to find $\hat{F}(\cdot)$, $\hat{A}(\cdot)$ and $\hat{P}(\cdot)$. The $\hat{F}(\cdot)$ is straightforward to obtain:

$$\hat{F} = \text{BER}(\alpha, \beta, \gamma), \quad (30)$$

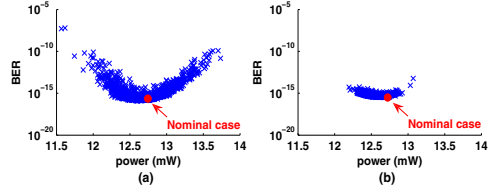


Fig. 5. Power and performance variation for 1000 die samples by Monte Carlo simulation: (a) without tuning and (b) with tuning.

where α is the vector indicating the LSB design for each tap. β and γ are vectors in R^n containing resolution and LSB size for each tap, and \bar{e} is the allowed BER upper bound. Note that the number of taps n is no longer a variable: by allowing $\beta_i = 0$, tap i is removed. Accordingly, we only need to specify n_{max} , a maximum number of taps to be considered ($n_{\text{max}} = 10$ in this paper).

The power P_{filter} (29) and area A_{filter} (27) of the pre-emphasis filter also need to be modified with the introduction of the DAC:

$$P_{\text{filter}} = \rho \sum_{i=1}^{n_{\text{max}}} D_i^T [2^{\beta_i-1}, \dots, 2^0] \cdot \gamma_i I_{\text{unit}} \cdot V_{dd} + (1 - \rho) f \cdot V_{dd}^2 \cdot C_{\text{para}}, \quad (31)$$

$$A_{\text{filter}}(\alpha, \beta, \gamma) = \sum_{i=1}^n 2^{\beta_i} \gamma_i A_{\text{unit}}^{\alpha_i}, \quad (32)$$

$$C_{\text{para}} = \sum_{i=1}^{n_{\text{max}}} 2^{\beta_i} \gamma_i C_{\text{unit}}^{\alpha_i}. \quad (33)$$

Note that vector D_i represents the digital control bits and P_{filter} becomes a distribution instead of a deterministic value because of the I_{unit} variation from V_{th} mismatch. The other calculations are kept the same and the total area and power can be obtained by (24) and (25), accordingly.

VI. PLL DESIGN IN HIGH-SPEED SERIAL LINK

An example of second-order PLL as shown in Fig. 6 comprises of several components: (1) the phase frequency detector, (2) the charge pump, (3) the loop filter, and (4) the voltage-controlled oscillator. Phase and frequency detector is used to detect phase and frequency difference and provides the UP/DN signal to the charge pump. The charge-pump circuit comprises of two switches driven by the UP and DN signal and injects the charge into or out of the loop filter capacitor (C_{CP}). The combination of charge-pump and C_{CP} is an integrator that generates the average voltage of UP (or DN) signal, V_{ctrl} , and adjusts the frequency of the subsequent oscillator circuit. In Fig. 6, a power-supply regulated ring oscillator is shown with the voltage-to-frequency gain K_{VCO} . The VCO output frequency is controlled by its supply voltage V_{ctrl} .

The performance of PLL is measured by its output clock jitter. Jitter can be expressed as $\sigma_{\Delta T} = (T/2\pi) \cdot \sigma_{\Delta\phi}$, where ω_0 is the clock frequency, $T = 2\pi/\omega_0$ is the clock period, and $\sigma_{\Delta\phi}$ is the phase jitter of the clock. The jitter mainly comes from the reference clock (N_{in}) and VCO (N_{VCO}), which can be expressed as [10]:

$$\text{Jitter} = \sigma_{\Delta T}^2 = \frac{8}{\omega_0^2} \int_0^\infty S_\phi(f) \sin^2(\pi f \Delta T) df, \quad (34)$$

$$S_\phi(f) = \frac{N_{in}}{f^2} \cdot |Hn_{in}(j2\pi f)|^2 + \frac{N_{VCO}}{f^2} \cdot |Hn_{VCO}(j2\pi f)|^2.$$

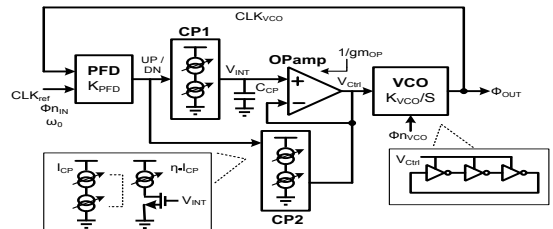


Fig. 6. Tunable and adaptive bandwidth PLL.[9]

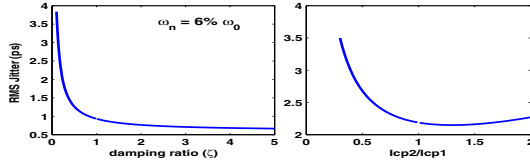


Fig. 7. Output jitter sensitivity to the (a) loop damping factor ζ and (b) charge pump current ratio I_{CP2}/I_{CP1} .

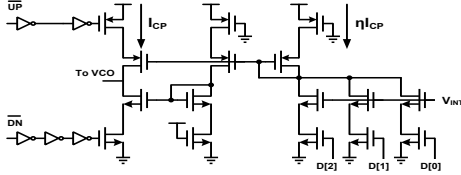


Fig. 8. Charge pump schematic[9].

Note that Hn_{in} and Hn_{VCO} are the noise transfer functions of the reference clock noise (N_{in}) and VCO noise (N_{VCO}) accordingly.

The jitter performance is a function of the PLL design parameters ω_n and ζ [10]. Fig. 7(a) shows an example of output root-mean-square (RMS) jitter with respect to the damping ratio (ζ) for a fixed $\omega_n = 0.06\omega_0$ [10]. Moreover, in the case of tunable PLL shown in Fig. 6, the natural frequency varies proportionally to $\sqrt{I_{CP1}}$ and the damping factor is proportional to $I_{CP2}/\sqrt{I_{CP1}}$ [9]. By finding an optimum value of the absolute value and relative ratio of I_{CP2} and I_{CP1} , we can minimize the PLL output jitter. We write the objective function of the design-time optimization as (1):

$$Prob(Jitter(\alpha) \leq \bar{f}), \quad (35)$$

where α is a vector which represents the number of unit cells used in the charge pumps. In other words, it represents the value of I_{CP1} and I_{CP2} . An example of the relation between output RMS jitter and the current ratio for the charge pumps (I_{CP2}/I_{CP1}) is shown in Fig. 7(b), with a fixed I_{CP1} .

For the design-time optimization, we want to minimize the output clock jitter, subject to power and area constraints. The design parameters are the charge pump currents I_{CP1} and I_{CP2} . The power consumption of the charge pump can be calculated by an approach similar to the one used in our first transmitter design example. Assume we use unit cells of type α_i ($1 \leq \alpha_i \leq m$) with unit current I_{α_i} , then the required number of cells q_i for the charge pump i can be determined by

$$q_i = \lceil \frac{I_{CPi}}{I_{\alpha_i}} \rceil. \quad (36)$$

As a result, the power consumed by the charge pump is:

$$P_{CP}(\alpha) = \sum_i \frac{1}{2} (2\pi\omega_0) \cdot (1 + \frac{1}{\eta_i}) \cdot q_i I_{\alpha_i} \cdot V_{dd}, \quad (37)$$

where η_i represents the current mirror ratio for the biasing circuit of the charge pump i . The area can also be approximated using the similar method and details can be found in our first example.

In the presence of process variation, the output RMS jitter varies for the same design because of the variations on I_{CP1} and I_{CP2} . To reduce the impact of process variation and improve the parametric yield, post-silicon tuning techniques can be applied. Fig. 8 shows a schematic of the charge pump circuit with digitally tuned elements placed in the biasing circuit. By applying a proper digital control signal D , the charge pump current ratio can be optimized to reduce the output jitter under the impact of process variation.

As discussed in Section III, we can change the objective function to the *Jitter parametric yield* as

$$Prob(Jitter(\eta, \alpha, \beta, \gamma) \leq \bar{f}), \quad (38)$$

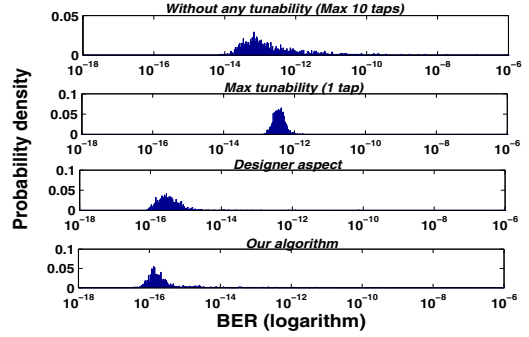


Fig. 9. BER distribution for four different designs.

where α is the vector indicating the LSB design for each tap in the tunable element. β and γ contain resolution and LSB size for each charge pump and η represents the biasing current ratio; \bar{f} is the allowed jitter upper bound. The power consumed by the charge pump can be re-written as

$$P_{CP} = \sum_i \frac{1}{2} (2\pi\omega_0) \cdot (1 + \frac{1}{\eta_i}) \cdot D_i^T [2^{\beta_i-1}, \dots, 2^0] \cdot \gamma_i I_{\alpha_i} \cdot V_{dd}.$$

Note that in this example, the tunable element is inserted in the biasing part with bias ratio η , which is considered as part of the design parameters α in (P1). When $\eta \ll 1$, only a small amount of current in the biasing circuit is required to generate I_{CP} . As a result, the power consumed and the area occupied by the digitally tuned element can be ignored. In this case, however, the LSB size in the charge pump current becomes $\frac{1}{\eta} \gamma I_{\alpha}$, which is increased when η is decreased. The effect of tuning is reduced and may not provide the desired yield. On the other hand, when $\eta \sim 1$, the tunability is maximized but the power and area consumed by the tunable element is also increased. Obviously, a good balance needs to be found through our proposed framework.

VII. EXPERIMENTAL RESULTS

We extract transistor model parameters by SPICE simulation in IBM 90nm technology and implement the proposed algorithm in MATLAB. All the experiments are run on a Windows server with Pentium IV 3.2GHz CPU and 2G RAM.

A. Transmitter Design

We compare our algorithm with three different methods: no-tunability design, maximum tunability design, and design heuristic from designer's perspective. The design heuristic is guided by the designers' experience [11], [12]: (1) total number of filter taps is iteratively determined by the channel response and the LMS algorithm. (2) assume that each tap of the filter has the same LSB size; (3) the LSB size is determined by considering the maximum and minimum filter coefficient. This design methodology serves as a heuristic for this joint optimization problem and essentially solves the problem in a reduced solution space. The no-tunability design sets the resolution to be 1 ($\beta_i = 1$) for all taps and maximizes the precision of a pre-set pre-emphasis filter. The maximum tunability design uses only one-tap filter ($n_{max} = 1$) to allow maximum adjustability. The no-tunability design and maximum tunability design also serve as the representative of maximum design-time effort and maximum post-silicon effort, respectively.

For fair comparison, the data rate for all the designs is set to be $5GHz$ and the threshold BER for yield $\bar{e} = 1.0 \times 10^{-15}$. In our experiments, we assume that the channel is a 30cm differential microstrip line on FR-4 substrates and that the receiver has ideal timing recovery. We also assume that V_{th} variation follows normal distribution.

We first present the BER distribution with 20% V_{th} variation based on 10K Monte Carlo runs in Fig. 9. The area is constrained to

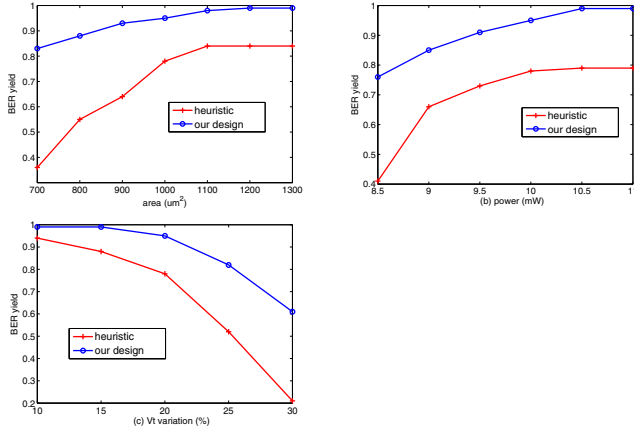


Fig. 10. Yield curves for our designs and design heuristic with respect to area (a), power (b) and V_{th} (c).

1000 μm^2 and the power is constrained to 10mW. First, for all the four methods, the distributions show strong non-symmetry and non-Gaussianity. This should be attributed to the non-linear relationship between the V_{th} and BER. Second, we can see that the ranges of BER vary for the four methods: the maximum and minimum tunability design gives the smallest and largest variations respectively, with the other two methods in between. This is in accordance with the intuition that more tunability corresponds to less variation. Third, we can see that our design gives the smallest mean BER while the minimum tunability design gives the largest mean BER. Moreover, compared with the design heuristic, our design optimizes the BER distribution with better mean and smaller variance. This verifies that our joint design-time and post-silicon optimization can significantly improve performance when compared with design-time or post-silicon only optimization, or the heuristic method.

Next, we quantitatively study how the yield from our design and design heuristic vary with respect to different area constraints for fixed power ($P = 10mW$) and 20% V_{th} variation. The yield is defined as the percentage of the chips meeting the BER as in (5). The results are presented in Fig. 10 (a), where for different area specs, our design always gives a larger yield than the design heuristic. When the area is limited to 700 μm^2 , we have a 47% yield improvement over the design heuristic. Note that when area constraint is larger than 1200 μm^2 , the yield does not improve because the design is dominated by the power constraint. This verifies our discussion that the power and area constraints are strongly coupled.

A similar study is conducted with respect to different power constraints for fixed area ($A = 1000\mu m^2$) and 20% V_{th} variation and to different V_{th} variation for fixed power ($P = 10mW$) and area ($A = 1000\mu m^2$). From Fig. 10 (b), we can see that for different power specifications, our design also gives better yield and better scalability than the design heuristic. When the power is limited to 8.5mW, we have a 35% yield improvement over the design heuristic. As for V_{th} variation, although it is not explicitly listed as a constraint and only appears in the power and area constraints, it affects the yield significantly. Our design improves the yield by 40% when compared with the design heuristic with 30% variation, as shown in Fig. 10 (c).

B. PLL Design

The same optimization framework is applied to a PLL design example and the result is provided in Fig. 11. We compare our algorithm with the design heuristic that has optimal I_{CP1} and I_{CP2} values through design time optimization and tunable elements in the biasing circuit consumes negligible power [10]. The reference clocks of the PLL for both designs are set to 700MHz. We assume that the V_{th} variation follows normal distribution. The yield is defined as the

percentage of the chips meeting the jitter requirement, as in (38). The experiment was conducted with respect to different power constraints for fixed area and 30% V_{th} variation, as shown in Fig. 11(a). From the figure we can see that for different power specs, our design provides better yield than the design heuristic and obtains up to 29% yield improvement. In Fig. 11(b), when the power is limited to 17mW, we have a 56% yield improvement over the design heuristic.

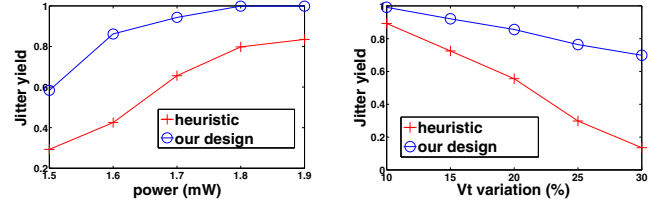


Fig. 11. Yield for our algorithms and design heuristic w.r.t power (a) and V_{th} (b) in the PLL.

VIII. CONCLUSIONS

Joint design time and post-silicon optimization for analog circuits has been an open problem in literature, given the complex nature of analog circuit modeling and optimization. In this paper we formulate a co-optimization problem for digitally tuned analog circuits to optimize the parametric yield, subject to power and area constraints. A general optimization framework combining the branch-and-bound algorithm and gradient ascent method is proposed. We demonstrate our framework with two examples in high-speed serial link, the transmitter design and the phase-locked-loop (PLL) design. Experimental results show that compared with the design heuristic from analog designers' perspective, joint design-time and post-silicon optimization can improve the yield by up to 47% for transmitter design and up to 56% for PLL design under the same area and power constraints.

REFERENCES

- [1] H. Darabi, S. Khorram, H.-M. Chien, M.-A. Pan, S. Wu, S. Moloudi, J. Leete, J. Rael, M. Syed, R. Lee, B. Ibrahim, M. Rofougaran, and A. Rofougaran, "A 2.4-ghz cmos transceiver for bluetooth," *Solid-State Circuits, IEEE Journal of*, vol. 36, pp. 2016–2024, Dec 2001.
- [2] H. Huang and E. K. F. Lee, "Design of low voltage cmos continuous-time filter with on-chip automatic tuning," *J. Solid-State Circuits*, 2001.
- [3] G. Miller, M. Timko, H.-S. Lee, E. Nestler, M. Mueck, and P. Ferguson, "Design and modeling of a 16-bit 1.5msps successive approximation adc with non-binary capacitor array," *Proc. Int. Great Lakes Symp. on VLSI*, 2003.
- [4] B. Murmann and B. Boser, "Digitally assisted analog integrated circuits," *Queue*, vol. 2, no. 1, pp. 64–71, 2004.
- [5] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge University Press, 2004.
- [6] Y. Tao, W. Bereza, R. Patel, S. Shumarayev, and T. Kwasniewski, "A signal integrity-based link performance simulation platform," *Proc. Custom Integrated Circuits Conference*, 2005.
- [7] S. Sen, V. Natarajan, R. Senguttuvan, and A. Chatterjee, "Pro-vizor: Process tunable virtually zero margin low power adaptive rf for wireless systems," in *Proc. Design Automation Conf.*, June 2008.
- [8] Y. Ye, F. Liu, S. Nassif, and Y. Cao, "Statistical modeling and simulation of threshold variation under dopant fluctuations and line-edge roughness," *Proc. Design Automation Conf.*, June 2008.
- [9] S. Sidiropoulos, D. Liu, J. Kim, G. Wei, and M. Horowitz, "Adaptive bandwidth dlls and pll's using regulated supply cmos buffers," pp. 124–127, 2000.
- [10] M. Mansuri and C.-K. K. Yang, "Jitter optimization based on phase-locked loop design parameters," *Solid-State Circuits, IEEE Journal of*, vol. 37, pp. 1375–1382, Nov 2002.
- [11] J. VIta, A. Marques, P. Azevedo, and J. Franca, *Design Considerations for a Retargetable 12b 200MHz CMOS Current-Steering DAC*. Springer US, 2003.
- [12] A. C. Y. Lin and M. J. Loinaz, "A serial data transmitter for multiple 10gb/s communication standards in 0.13 μm cmos," in *Int. Solid State Circuits Conf.*, 2008.