

A Multiphase Decision Model for System Reliability Growth With Latent Failures

Tongdan Jin, *Member, IEEE*, Ying Yu, and Hong-Zhong Huang

Abstract—Reliability growth testing becomes difficult to implement as the product development cycle continues to shrink. As a result, the new design is prone to latent failures due to design immaturity and uncertain operating condition. Reliability growth planning emerged as a new methodology to drive the reliability across the product lifetime. We propose a multiphase reliability growth model that sequentially determines and implements corrective actions (CAs) against surfaced and latent failure modes. Such a holistic approach enables the manufacturer to attain the reliability goal while ensuring the product time to market. We devise a CA effectiveness function to assess the tradeoff between the failure removal rate and the required resources. Rosen's gradient projection algorithm is used to determine the optimal resource allocation in each phase. The applicability and performance of the reliability growth model are demonstrated on a fleet of semiconductor testing equipment.

Index Terms—Capital equipment, corrective action (CA) effectiveness, latent failure, power law model, reliability growth planning (RGP).

N_i	Number of failures for the failure mode i between $[0, t_c]$.
t_{in}	n th failure arrival time for failure mode i for $n = 1, 2, \dots, N_i$.
T_c	Time interval between t_1 and t_c .
T	Time interval between t_c and t .
L	Set of latent failure modes that occurred in T_c .
$\hat{\Gamma}(t)$	Cumulative failure intensity estimate for latent failure modes.
k_c	Number of latent failure modes in T_c .
$h(x)$	CA effectiveness function.
$g(x)$	CA ineffectiveness function, $= 1 - h(x)$.
b, c	Parameters for $h(x)$.
$f(\mathbf{x}; t), f(\mathbf{x})$	Objective function with decision vector \mathbf{x} .
C	Total CA budget for a particular decision phase.
\mathbf{P}	Rosen's gradient projection matrix.

NOMENCLATURE

m	Number of surfaced failure modes by time t_c .
k	Number of latent failure modes that will occur between t_c and t .
t_1, t_c , and t	Previous time, current time, and future time, respectively.
$\mu_i(t)$	Failure intensity for surfaced failure mode i for $i = 1, 2, \dots, m$.
$\gamma_j(t)$	Failure intensity for latent failure mode j for $j = 1, 2, \dots, k$.
$\hat{\mu}_i(t)$	Estimate of $\mu_i(t)$.
$\hat{\gamma}_j(t)$	Estimate of $\gamma_j(t)$.
$\hat{\mu}_s(t t_c)$	Estimate of system failure intensity at t .
$\hat{\mu}_{s,CA}(\mathbf{x}; t)$	Estimate of system failure intensity after CA.
x_i	CA resource against failure mode i , a decision variable.
α, β	Parameters of the Crow/AMSAA model.
$\hat{\alpha}, \hat{\beta}$	Estimates for α and β , respectively.

Manuscript received February 26, 2010; revised August 8, 2011 and March 21, 2012; accepted July 13, 2012. Date of publication April 18, 2013; date of current version June 12, 2013. This paper was recommended by Associate Editor H. Pham.

T. Jin is with the Ingram School of Engineering, Texas State University—San Marcos, San Marcos, TX 78666 USA (e-mail: tj17@txstate.edu).

Y. Yu is with the Department of Automation, Shanghai University, Shanghai 20000, China (e-mail: foolcatonline@hotmail.com).

H.-Z. Huang is with the School of Mechatronics Engineering, University of Electronic Science and Technology of China, Chengdu 61000, China (e-mail: hzhuang@uestc.edu.cn).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TSMCA.2012.2217318

I. INTRODUCTION

RELIABILITY growth test (RGT) aims to identify the design weakness, remove the critical failure modes, and enhance the system performance prior to the volume production. Depending on the occurrence time, all types of failure modes can be classified into two categories: surfaced failure mode and latent failure mode. In this paper, a surfaced failure mode is defined as a failure mode which occurs during the in-house RGT process. A latent failure mode is also called a dormant failure mode. It usually occurs post the system installation or if the in-house testing time could be extended. Issues related to latent failures have been frequently reported in electronics industry [1], [2], robotic and mobile vehicles [3], and computer servers [4].

Both the quantity and the incident rate of latent failures are highly stochastic in nature. One particular reason that triggers a latent failure is electrical static discharge (ESD). ESD often damages the electronic devices or components during the system manufacturing, handling, and installation, but the failure symptom does not show up until the field operation [5]. Other issues such as software bugs, design weakness, or improper usage could also trigger or induce latent failures [6]. Latent failures could be mitigated or avoided if the system is built or equipped with health prognostics capability or condition-based monitoring tools. Readers are referred to [7] for detailed discussions on this technology. A latent failure mode, once it becomes a dominant failure mode, may create a large financial pressure on the manufacturer, not to mention the loss of customer goodwill. Such expenses include escalated warranty costs, excessive spare parts inventory, and increased

maintenance labor. Hence, it is highly desirable to understand the nature of latent failures so that countermeasures can be adopted prior to the plague of the issue.

RGT is also known as test-analyze-and-fix. It can be applied to a new system design if sufficient in-house testing time is available. The idea of RGT can date back to Duane [8] when he attempted to identify the relationship between the reliability growth rate and the testing time in aircraft components. Crow [9] found that reliability growth rate can be approximated by a nonhomogenous Poisson process (NHPP). His method became the well-known Crow/Army Materiel Systems Analysis Activity (AMSAA) model which is also called the power law model. Since then, many studies pertaining to the RGT methodology have been reported [10]–[15]. For instance, Xie and Zhao [10] developed a graphical tool to predict the reliability growth based on the Duane model. Campbell [11] proposed an optimization model to allocate subsystem test time for maximizing the system reliability. Coit [12] generalized Campbell's model by adding the test budget as a design constraint. Benski and Cabau [13] used the design of experiment to determine the optimal testing parameters in RGT programs. Krasich *et al.* [14] and Krasich [15] proposed accelerated RGT procedures to reduce the testing time subject to reliability constraints. While the aforementioned works focus on hardware testing, the studies in [16]–[18] extended the RGT concept to software design, debugging, and performance verification. In general, these models are quite effective to drive the reliability of new products when sufficient in-house testing time and relevant resources are available. Throughout this paper, product and system are used interchangeably.

The traditional RGT process becomes difficult to implement in industries where the new design is pushed under a fast time-to-market pressure. Examples include semiconductor equipment, consumer electronics, wind and solar generation systems, electric vehicles, and medical devices. For instance, automatic test equipment (ATE) is a high-end electronics machine commonly used in wafer testing industry. ATE makers constantly redesign and upgrade the equipment to meet the new chip performance requirement governed by Moore's law. In today's fast-paced yet distributed business environment, ATE makers cannot fully rely on the in-house RGT process to attain the reliability goal. In addition, complex systems such as ATE are often designed in modularity to facilitate the maintenance and repair. The reliability testing becomes more challenging when different types of modules have a different development timeline. As such, it is almost impossible to assemble all types of modules to perform a system-level test.

To maintain the competitive edge, a new design must be released to the market in a timely manner. Meanwhile, costs related to design, manufacturing, and testing must be minimized yet without compromising the reliability performance [19]. To address these problems, it is imperative to develop a new reliability management program that is capable of meeting the product delivery deadline and ensuring the reliability goal. In the last decade, reliability growth planning (RGP) has emerged as a new methodology to resolve these challenging issues [20]–[22]. RGP differs from RGT in that it drives the reliability across design, manufacturing, and field operation. This new

concept allows the system manufacturer to implement corrective actions (CAs) and, if necessary, reallocate the CA resources prior and post the product shipment. As such, a new design can be released to the market in a timely manner, and the reliability is improved through lifetime commitment.

Under the RGP scheme, this paper proposes a multiphase reliability growth model that guides the manufacturer to attain the reliability goal through a series of CA initiatives. While a large body of literature addressing latent failures is available, these studies usually focus on the prediction and analysis of latent failures. The objective of this paper is not only to predict the latent failure but also to demonstrate how the CA resources need to be redistributed given the occurrence of new failure modes. We synthesize the CA effectiveness, the reliability growth, and the failure prediction into a unified optimization framework. The contributions to the reliability community are twofold. First, we propose a CA effectiveness function to characterize the failure reduction rate per unit amount of CA budget. Such quantification allows the decision maker to identify and attack the failure modes that consume the smallest budget yet achieve the largest reliability growth. Second, we formulate a multiphase CA optimization model to mitigate, if not eliminate, all critical latent failures in a sequential manner. This allows the manufacturer to address the surfaced and the emerging failure modes at the same time.

The remainder of this paper is presented as follows. Section II briefly introduces the multiphase RGP concept. Section III reviews the latent failure prediction model. Section IV proposes an analytical model to characterize the CA effectiveness. In Section V, an optimal decision model synthesizing the CA cost with the latent failure modes is formulated. In Section VI, the proposed model is demonstrated on field ATE systems. Section VII concludes this paper with some remarks on the future research.

II. CONCEPT OF MULTIPHASE RELIABILITY GROWTH

Recently, there has been a small but growing stream of RGP literature reported from industry and academia. For example, Smith [20] applied the RGP concept to calculate the CA-based reliability improvement cost for a fleet of systems. Ellner and Hall [21] proposed an RGP model to estimate system reliability growth taking into account latent failures. Jin and Wang [22] formulated a multicriteria programming model to maximize field system reliability through optimal CA decisions. These studies show that RGP is an effective approach to drive the product reliability if extended in-house testing is infeasible in early design and prototyping phase.

RGP is a lifetime commitment that continuously improves the system reliability through in-depth failure analysis and rigorous CA programs across the product lifecycle. Two types of CA are generally applied to repairable systems: retrofit and engineering change order (ECO) [21], [22]. Retrofit uses spare modules (i.e., line replaceable units) to proactively replace in-service modules that will fail due to a known failure mode. ECO is a countermeasure often implemented in the repair center to eliminate critical failure modes when modules are returned from field. In general, retrofit is more costly than ECO because

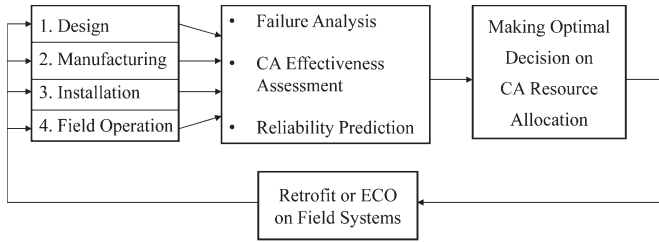


Fig. 1. Multiphase RGP process.

it requires dedicated personnel, logistics, and spare parts for performing on-site replacement tasks.

This paper aims to extend the preliminary findings [22] to a multiphase reliability growth environment. The study is motivated to address the stochastic behavior of latent failures and further mitigate their impact on product reliability. Latent failures have been studied in the literature for decades, but they have not been adequately addressed in terms of the prediction and elimination under RGP scheme. CA could be simply focused on surfaced failure modes if latent failure modes are negligible. However, CA resources must be redistributed if a latent failure mode becomes a dominant issue. The proposed multiphase decision model aims to fill this gap by redistributing the CA resources to emerging failure modes if needed.

The flowchart in Fig. 1 explains the technical realization of the proposed multiphase RGP model. Failure data observed from system design, manufacturing, and field operation are collected, and their root causes are analyzed. Optimal CA decisions are determined with the goal to minimize the failure intensity or equivalently maximize the system reliability. Following the CA implementation, failure intensities of surfaced and latent failure modes are monitored and compared to their anticipated value. CA resources are redirected to latent failure modes which might become dominant issues. This process is repeated over the planning horizon until the system reliability goal is achieved.

III. PREDICTION OF LATENT FAILURES

Latent failure modes are often “embedded” in the system, and the times of occurrence are highly stochastic in nature. Therefore, it is usually difficult to predict the failure intensity. In [23] and [24], Markov models have been used to predict the aircraft component reliability with the consideration of latent failures. The Markov method is quite effective when the occurrence rate of latent failures is constant. In [25], a more general prediction model considering time-varying intensities is proposed. This model is briefly reviewed as it will be used to construct the multiphase RGP model. Assuming that a system has operated through t_c , then, the system failure intensity at t for $t > t_c$ can be forecasted as

$$\hat{\mu}_s(t|t_c) = \sum_{i=1}^m \hat{\mu}_i(t) + \sum_{j=1}^k \hat{\gamma}_j(t). \quad (1)$$

Equation (1) comprises two types of failures, i.e., surfaced failure modes and latent failure modes. m is the number of surfaced failure modes observed by t_c , and k is the number

of latent failure modes expected to occur between $[t_c, t]$. Both $\hat{\mu}_i(t)$ and $\hat{\gamma}_j(t)$ represent the failure intensity estimates for surfaced and latent failure modes, respectively. $\hat{\mu}_i(t)$ is the estimate for $\mu_i(t)$ which is the true yet unknown failure intensity; so is $\hat{\gamma}_j(t)$ for $\gamma_j(t)$. In practice, $\hat{\mu}_i(t)$ can be estimated based on failure data between $[0, t_c]$. Equation (1) is derived assuming that all failure modes are mutually independent. If correlations between two failure modes are relatively small, this model is still valid.

Without loss of generality, the NHPP model [9] is used to estimate $\mu_i(t)$ in this study. It is worth mentioning that (1) represents a more general prediction method and it can accommodate other trend functions such as bounded intensity process (BIP) model [26], [27]. The estimation process for $\hat{\gamma}_j(t)$ will be further discussed in Section III-B.

A. Prediction of Surfaced Failure Modes

The Crow/AMSAA test is probably the most widely used tool for assessing the reliability growth trend based on surfaced failure mode information. The essence of the Crow/AMSAA test is to determine whether a failure process is homogenous Poisson process (HPP) or NHPP. The underlying assumption is that the failure intensity function $\mu(t) = \alpha\beta t^{\beta-1}$ is adequate to capture the failure incidence behavior. When $\beta = 1$ and $\mu(t) = \alpha$, the process simply is an HPP. For $\beta > 1$, it implies that the failure intensity is increasing. If $\beta < 1$, $\mu(t)$ is a decreasing function with less failures occurring in the same length of interval. According to [9], the maximum likelihood estimates for β_i and α_i are

$$\hat{\beta}_i = \frac{N_i}{\sum_{n=1}^{N_i} \ln\left(\frac{t_c}{t_{in}}\right)} \quad \hat{\alpha}_i = \frac{N_i}{t_c^{\hat{\beta}_i}} \quad (2)$$

where N_i is the number of failures observed by t_c pertaining to failure mode i . Assuming that the observation starts at $t = 0$, then, t_{in} is the n th failure arrival time. More discussions about the trend test and analysis are available in [9] and [28]. When the Crow/AMSAA model is appropriate for predicting the surfaced failure mode intensity, (1) can be rewritten as

$$\hat{\mu}_s(t|t_c) = \sum_{i=1}^m \hat{\alpha}_i \hat{\beta}_i t^{\hat{\beta}_i-1} + \sum_{j=1}^k \hat{\gamma}_j(t) \quad (3)$$

where $\hat{\alpha}_i$ and $\hat{\beta}_i$ are the parameters for the i th surfaced failure mode for $i = 1, 2, \dots, m$. In some situations, the failure intensity may eventually level off. Then, it is more appropriate to use the BIP model to predict the failure intensity [26], [27].

B. Prediction of Latent Failure Modes

The second summation in (3) represents the failure intensity of latent failure modes that are expected to occur between $[t_c, t]$. Notice that k is the expected number of latent failure modes. The value of k can be appropriately estimated based on existing surfaced failure modes, and the result is given as follows:

$$k \cong \left\lfloor \frac{k_c(t - t_c)}{t_c - t_1} \right\rfloor = \left\lfloor \frac{k_c T}{T_c} \right\rfloor. \quad (4)$$

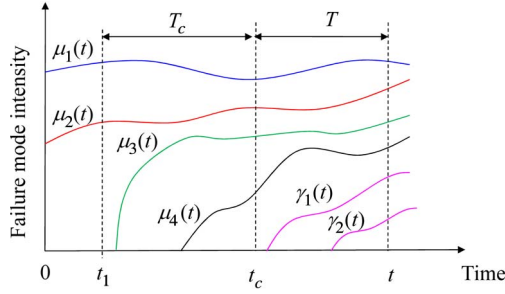


Fig. 2. Surfaced and latent failure mode intensities.

Both T and T_c are defined in Fig. 2. Note that $\lfloor \bullet \rfloor$ represents the integer part and k_c is the number of latent failure modes observed in T_c or between $[t_1, t_c]$. Similarly, $\hat{\gamma}_j(t)$ is the failure intensity for the j th latent failure mode with $j = 1, 2, \dots, k$. It is often difficult, if not possible, to estimate the failure intensity of individual $\gamma_j(t)$. However, the aggregate failure intensity, denoted as $\hat{\Gamma}(t) = \sum_{j=1}^{k_c} \hat{\gamma}_j(t)$, can be predicted by [25]

$$\hat{\Gamma}(t) \cong \frac{T}{T_c} \sum_{j=1}^{k_c} \hat{\mu}_j(t - T_c) = \frac{T}{T_c} \sum_{j \in L} \hat{\mu}_j(t - T_c) \quad (5)$$

where $\hat{\mu}_j(t - T_c)$ is the j th estimated latent failure mode expected to occur during T and L is the set of the latent failure modes occurred during T_c . For example, in Fig. 2, $\mu_3(t)$ and $\mu_4(t)$ are eligible as the latent failures to forecast $\hat{\Gamma}(t)$, but $\mu_1(t)$ and $\mu_2(t)$ are not because they occurred prior to t_1 . Thus, the set is $L = \{3, 4\}$. It is preferable to choose t_1 such that T_c is equal or close to T because new failure modes that emerged in T_c are more informative to forecast the latent failures in T . For detailed discussion on (4) and (5), readers are referred to [25]. By substituting (5) into (3), a system failure intensity function incorporating potential latent failures is obtained as

$$\hat{\mu}_s(t|t_c) = \sum_{i=1}^m \hat{\alpha}_i \hat{\beta}_i t^{\hat{\beta}_i - 1} + \frac{T}{T_c} \sum_{j \in L} \hat{\mu}_j(t - T_c). \quad (6)$$

Equation (6) is updated iteratively over the planning horizon. As the time evolves from t_c to t , latent failure modes occurring in $[t_c, t]$ are classified as surfaced failure modes based on which new latent failures in the next phase are predicted.

IV. CA EFFECTIVENESS FUNCTION

The CA effectiveness function aims to link the failure reduction rate of a particular failure mode with the amount of CA resources (e.g., money) required. The maximum CA effectiveness is one if that particular failure mode is completely eliminated from field systems. The minimum effectiveness is zero if no CA is applied. The effectiveness function plays a substantial role in determining the best resource allocation policy. To the best of our knowledge, the study [29] makes the first attempt to establish the relationship between the CA

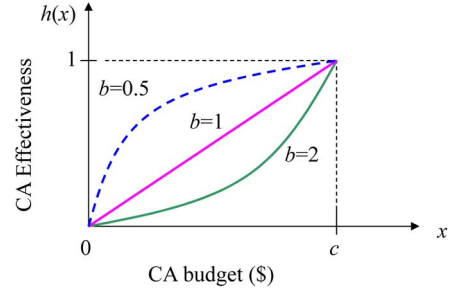


Fig. 3. Various shapes of CA effectiveness functions.

effectiveness and the associated resources. The model is restated as follows:

$$h(x) = \left(\frac{x}{c}\right)^b. \quad (7)$$

In this model, x represents the amount of the CA budget allocated to a particular failure mode. Both b and c are model parameters which can be estimated from historical CA data, or they can be extrapolated from predecessor products. The value of c actually is equal to the retrofit cost under the assumption that all field systems receive retrofit service. This can be easily justified from the fact that when $x = c$, $h(x) = 1$.

By changing b , three types of effectiveness functions are available as shown in Fig. 3: linear, quadratic, and rational. If $b = 1$, $h(x)$ simply becomes a linear function with the implication that the failure removal rate is proportional to the CA budget. If $b > 1$, $h(x)$ turns out to be a power function. For $b < 1$, $h(x)$ becomes a rational function. This can be applied to situations where the CA effectiveness decreases once the budget reaches a certain level. Typical examples include software reliability growth and process yield improvement.

V. OPTIMIZATION FORMULATION

A. Reliability Prediction Considering CAs

In determining the CA resources, priorities are often given to those failure modes showing high failure intensity rates. To obtain a generalized system failure intensity estimate, it is assumed that CA could be applied to all surfaced failure modes regardless of their intensity. By combining (6) and (7), the system failure intensity function upon the CA is given

$$\begin{aligned} \hat{\mu}_{s,CA}(\mathbf{x}; t) &= \sum_{i=1}^m (1 - h_i(x_i)) \hat{\mu}_i(t) + \frac{T}{T_c} \sum_{j \in L} \hat{\mu}_j(t - T_c) \\ &= \sum_{i=1}^m g_i(x_i) \hat{\mu}_i(t) + \frac{T}{T_c} \sum_{j \in L} \hat{\mu}_j(t - T_c). \end{aligned} \quad (8)$$

Equation (8) brings the CA budget x_i effects into the future system failure intensity. Here, $g_i(x_i) = 1 - h(x_i)$ is the CA ineffectiveness function for the i th surfaced failure mode with $\mathbf{x} = [x_1, x_2, \dots, x_m]$. An assumption behind (8) is that any ongoing CA neither induces nor eliminates latent failure modes. Uncertainty is often involved in reliability estimation

and prediction. The mean and the variance of the estimate are commonly used to characterize the uncertainty. The mean value and its variance of $\hat{u}_{s,CA}(\mathbf{x}; t)$ can be estimated by

$$E[\hat{\mu}_{s,CA}(\mathbf{x}; t)] = \sum_{i=1}^m g_i(x_i) E[\hat{\mu}_i(t)] + \frac{T}{T_c} \sum_{j \in L} E[\hat{\mu}_j(t - T_c)] \quad (9)$$

$$\text{var}(\hat{\mu}_{s,CA}(\mathbf{x}; t)) = \sum_{i=1}^m (g_i(x_i))^2 \text{var}(\hat{\mu}_i(t)) + \left(\frac{T}{T_c}\right)^2 \sum_{j \in L} \text{var}(\hat{\mu}_j(t - T_c)). \quad (10)$$

Both equations are derived assuming that all failure modes are mutually independent or the correlation is small to be ignored. At a given time instance, the failure intensity of individual failure modes is a random variable. The mean and the variance of $\hat{u}_{s,CA}(\mathbf{x}; t)$ are the sum of the mean and the variance of individual failure modes, respectively. If m is reasonably large, $\hat{u}_{s,CA}(\mathbf{x}; t)$ tends to be normally distributed based on the central limit theorem (CLT). In probability theory, CLT states the conditions under which the sum of a sufficiently large number of independent random variables, each with finite mean and variance, tends to be normally distributed [30]. At given t , the CLT naturally leads to the following result:

$$\hat{\mu}_{s,CA}(\mathbf{x}; t) \sim \text{Normal}(E[(\hat{\mu}_{s,CA}(\mathbf{x}; t)], \text{var}(\hat{\mu}_{s,CA}(\mathbf{x}; t))). \quad (11)$$

Based on (11), the decision maker can allocate the CA resources in a way that the mean and the variance of $\hat{u}_{s,CA}(\mathbf{x}; t)$ are minimized or, equivalently, the system reliability is maximized. In the following, we formulate an optimization model to guide the implementation of the CA process.

B. Optimization Formulation

RGP aims to find the best way to implement the CA programs so as to attain the design goal. This is equivalent to minimizing the system failure intensity by appropriately allocating CA resources against critical failure modes. Since most failure modes behave randomly, the uncertainty of the failure intensity should be incorporated into the optimization model. These criteria can be synthesized into a unified decision model by minimizing the upper bound of $\hat{u}_{s,CA}(\mathbf{x}; t)$. Now, the optimization model, denoted as Problem P1, is formulated as (12)–(14), shown at the bottom of the page, where $Z_{1-\theta}$ is the standard normal value

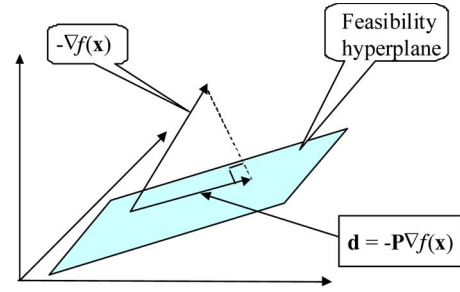


Fig. 4. Gradient projection method.

at $(1 - \theta) \times 100\%$ confidence level. The decision variable x_i represents the CA budget allocated against failure mode i for $i = 1, 2, \dots, m$. Equation (13) is a linear constraint limiting the maximum budget at the current decision phase. Problem P1 would be particularly beneficial to the mitigation of the reliability uncertainty if sources or qualities of the estimated model parameters differ appreciably within the system. As pointed by Coit [12], a decision strategy which ignores the uncertainty may put unwarranted resources on one failure mode with potential improvement. This resource allocation is promising, but it could also be risky. The risk could be mitigated if a more conservative plan can be implemented which assures the resource allocation to other failure modes. Obviously, $f(\mathbf{x}; t)$ is formulated to incorporate the uncertainty of various failure modes.

C. Optimization Algorithm

Problem P1 can be solved by successively projecting the gradient of the objective function onto a hyperplane constituted by possible solutions. The hyperplane is constructed by the cost equality in (13). This method was proposed by Rosen [31], and the concept of the gradient projection is explained in Fig. 4.

A projection matrix \mathbf{P} is used as a premultiplier to project the gradient onto the feasible hyperplane. The matrix \mathbf{P} for a single linear constraint in (13) can be readily determined. For our problem, \mathbf{P} is an m -by- m matrix which is defined by the following equation:

$$\mathbf{P} = \frac{1}{m} \begin{bmatrix} m-1 & -1 & \dots & -1 \\ -1 & m-1 & & \vdots \\ \vdots & & \ddots & -1 \\ -1 & \dots & -1 & m-1 \end{bmatrix}. \quad (15)$$

Since $x_i \geq 0$ is required for all i during the searching process, \mathbf{P} is dynamically updated to avoid negative x_i in each iteration. If negative x_i occurs, they are set to zero and the cost associated with those negative variables is redistributed to all

Problem P1 :

$$\text{Min} : f(\mathbf{x}; t) = E[(\hat{\mu}_{s,CA}(\mathbf{x}; t)] + Z_{1-\theta} (\text{var}(\hat{\mu}_{s,CA}(\mathbf{x}; t)))^{\frac{1}{2}} \quad (12)$$

Subject to :

$$\sum_{i=1}^m x_i \leq C, \quad (13)$$

$$x_i \geq 0 \text{ for } i = 1, 2, \dots, m \quad (14)$$

TABLE I
SUMMARY OF FAILURE MODE DISTRIBUTION

Category	i	Failure Mode	First Failure Time (days)	Cumulative Failures by day 350
Surfaced Failure Modes	1	Open Diode	7	17
	2	Power Supply	14	7
	3	EEPROM	14	2
	4	Cold Solder	21	1
	5	NFF	84	6
	6	Flux Contam	84	1
Latent Failure Modes	7	SMC Limit Table	105	1
	8	Capacitor	161	1
	9	PPMU	168	2
	10	Missing Solder	168	1
	11	Mfg Defects	210	1
	12	Bad ASIC	231	1
	13	Fuse	266	1
	14	Open Trace	287	1
	15	Op-Amp	315	2
	16	Timing Generator	315	1
	17	Solder Short	343	1
Total failures				47

positive x_i . The negative gradient projected onto the feasible region is given in (16) along with the gradient vector in (17)

$$\mathbf{d} = -\mathbf{P}\nabla f(\mathbf{x}) \quad (16)$$

$$\nabla f(\mathbf{x}) = \left[\frac{\partial f(\mathbf{x})}{\partial x_1} \quad \frac{\partial f(\mathbf{x})}{\partial x_2} \quad \frac{\partial f(\mathbf{x})}{\partial x_3} \quad \cdots \quad \frac{\partial f(\mathbf{x})}{\partial x_m} \right]^T \quad (17)$$

where $\mathbf{x} \in R^m$ and $\mathbf{d} \in R^m$. Successive iterations are continued until the optimality conditions are satisfied. According to [32], when $\mathbf{P}\nabla f(\mathbf{x}) = 0$, the Kuhn–Tucker optimality conditions have been met, and the optimal solution is found. When $\mathbf{P}\nabla f(\mathbf{x}) < \varepsilon$, it is assumed that the current solution is the optimal solution. In practice, a small threshold value ε in the range of 10^{-6} – 10^{-10} should be selected to determine the termination of iteration. Besides the Rosen’s algorithm, heuristic methods and genetic algorithm can also be considered for solving this type of multiphase RGP problems [19], [33], [34].

VI. CASE STUDIES

ATE is a complex electromechanical system widely used to test wafers at the back end of the semiconductor manufacturing process. A high-end ATE system typically costs \$2–3 million, including the purchasing, the training and operation, and the maintenance and upgrading during its useful lifetime. Driven by the short-cycled semiconductor market, ATE makers usually ship the new equipment to the market once the functional performance meets the customer requirement while the system reliability could be still low. ATE makers continue to improve the reliability performance by implementing rigorous CA programs. This process may last two or three years until field systems reach the reliability target.

A. Preliminary Data Analysis

The data set in Table I, adopted from [25], is used to demonstrate the application of the multiphase RGP model. These failure data were collected from 24 systems in 350 days after the initial installation. All field failures are classified into 17

failure modes. Notice that the first failure time differs among different failure modes. The first failure time is the time when a particular failure mode began to occur in the field systems. For example, the first failure time for open diode occurs in day 7 after the system installation. By the end of 350 days, there are a total of 17 open diode failures reported from 24 field systems.

An interesting observation is that the top six failure modes began to occur within 91 days or three months and the rest of eleven failure modes occurred after day 92. In this example, if a new failure mode occurred prior to 91 days, it is treated as a surfaced failure mode. This criterion is determined based on the fact that the ATE maker treats failures that occur during the first three months as the initial installation failures which create a major impact on the customer satisfaction. If the annual operating time per system is 8760 h, the mean time between failures (MTBF) for the fleet systems would be 6000 h based on surfaced failures. The actual MTBF is reduced by 27% with only 4400 h after taking into account all latent failure modes. This example indicates that latent failures need to be appropriately addressed when planning and executing reliability growth projects.

B. Multiphase RGP Implementation

In this section, the multiphase RGP optimization model in Section V-B is applied to drive the reliability growth of ATE systems. The entire RGP program consists of three interrelated phases: Phase 1 from day 1 to 91, Phase 2 from day 91 to 210, and Phase 3 from day 210 to 350. These phases are determined based on the ATE industry practice where three to four months are often adopted as a project implementation period. In Phase 1, system reliability is analyzed using the preliminary failure data. In Phase 2, an optimal CA decision is made based on the failure data from the previous phase. Phase 3 evaluates the ongoing CA effectiveness and decides whether resources need to be redistributed to emerging latent failure modes.

Each system is assumed to operate 24 h a day and seven days a week. This is the typical production environment in semiconductor industry. When an ATE system fails, the defective module is immediately replaced by a spare part, and the down time is small and can be ignored in our analysis. The failure modes in Table I along with their interarrival times are utilized to illustrate the multiphase RGP optimization program.

Phase 1—Reliability Prediction: To implement the RGP model, the first step is to analyze and predict the failure intensity based on the failure data in Phase 1. Six surfaced failure modes have been observed during the first 91 days. To estimate the Crow/AMSAA failure intensity, the interarrival times between two consecutive failures are required, and the details are presented in Table II.

Now, (2) can be used to estimate α and β for these surfaced failure modes, and the results are presented in Table III. In this phase, $t_1 = 0$; hence, $t_c = T_c$ and both are equal to 52 416 h. The prediction is made for the end of Phase 2, that is, $t = 120 960$ h and $T = 68 544$ h. All these are estimated based on the fleet cumulative operating hours. If the failure quality of a surfaced failure mode is one in Phase 1, it is assumed that the failure follows the HPP with $\beta = 1$. The failure intensity

TABLE II
FAILURE INTERARRIVAL TIMES IN PHASE 1

<i>i</i>	Days	7	14	15	21	84	85	87	89
1	Open Diode	1			1	1	1	1	1
2	Power Supply		1						
3	EEPROM		1	1					
4	Cold Solder			1					
5	NFF					1			
6	Flux Contam					1			

TABLE III
CROW/AMSAA RELIABILITY FORECASTING FOR PHASE 2

<i>i</i>	Failure Mode	$\hat{\alpha}_i$	$\hat{\beta}_i$	$E[\hat{\mu}_i(t)]$	$var(\hat{\mu}_i(t))$
1	Open Diode	1.29E-6	1.413	2.28E-4	5.2174E-10
2	Power Supply	1.91E-5	1.00	1.91E-5	3.6398E-12
3	EEPROM	5.40E-3	0.544	1.42E-5	2.0126E-12
4	Cold Solder	1.91E-5	1.00	1.91E-5	3.6398E-12
5	NFF	1.91E-5	1.00	1.91E-5	3.6398E-12
6	Flux Contam	1.91E-5	1.00	1.91E-5	3.6398E-12
7	Latent Failures	2.39E-4	1.021	3.07E-4	9.4433E-10

TABLE IV
OPTIMAL CA BUDGET ALLOCATION IN PHASE 2 WITH $\theta = 0.05$

<i>i</i>	Failure Mode	c_i (\$)	b_i	x_i (\$)
1	Open Diode	430,000	1	412,790
2	Power Supply	150,000	1	0
3	EEPROM	250,000	1	0
4	Cold Solder	75,000	1	19,510
5	NFF	370,000	1	0
6	Flux Contamination	45,000	1	27,700
7	Latent Failures in Phase 2	N/A	N/A	N/A

for potential latent failures in Phase 2 is also forecasted using (4) and (5), and the result is presented in Table III. For the illustration purpose, the standard deviation for individual failure modes is simply assumed as 10% of its mean value. Methods to compute the exact variance of $\hat{\mu}_i(t)$ are available in [35].

Phase 2—Decision Making on CAs: Phase 2 concentrates on resource allocation and CA implementation. To solve the optimization problem in Problem P1, we need to specify the values of b and c in the CA effectiveness model in (7). The value of c is equal to the retrofit cost, and it is relatively easy to estimate for individual failure modes. Estimating b is more involved, particularly if historical data are not available. In this case, we let $b = 1$ for all the failure modes, meaning that the effectiveness is proportional to the amount of the allocated CA budget. This assumption is similar to the Bayesian inference where the uniform distribution is adopted as the prior brief if historical data are not available. The total CA budget for Phase 2 is $C = \$460\,000$. After substituting the information from Tables III and IV into Problem P1, we apply the Rosen’s projection algorithm to search the optimal solution, and the result is listed in the last column of Table IV.

The solution suggests that the CA effort should be focused on the open diode issue. This is understandable because open diode is the dominant failure mode in Phase 1 showing an increasing trend with $\beta = 1.413$. On the other hand, the solution suggests that certain amounts of resources should be allocated to cold solder and flux contamination. Both failure intensities are relatively lower compared to others. Since their CA cost is also relatively lower compared to no fault found (NFF) and power supply, the decision is made such that CA

TABLE V
FAILURE INTERARRIVAL TIMES IN PHASE 2

<i>i</i>	Days	105	161	162	168	209	210
1	Open Diode	1				1	1
2	Power Supply	1	1	1			
3	EEPROM						
4	Cold Solder						
5	NFF		1		1		
6	Flux Contam						
7	SMC Limit Table	1					
8	Capacitor		1				
9	PPMU				1		
10	Missing Solder				1		
11	Mfg Defects						1

TABLE VI
CROW/AMSAA RELIABILITY FORECASTING FOR PHASE 3

<i>i</i>	Failure Mode	$\hat{\alpha}_i$	$\hat{\beta}_i$	$E[\hat{\mu}_i(t)]$	$var(\hat{\mu}_i(t))$
1	Open Diode	2.30E-4	0.903	6.40E-5	4.09E-11
2	Power Supply	2.66E-5	1.02	3.40E-5	1.16E-11
3	EEPROM	2.51E-2	0.374	4.49E-6	2.02E-13
4	Cold Solder	8.27E-6	1.00	8.27E-6	6.83E-13
5	NFF	4.22E-11	2.14	9.46E-5	8.94E-11
6	Flux Contam	8.27E-6	1.00	8.27E-6	6.83E-13
7	SMC Limit Table	8.27E-6	1.00	8.27E-6	6.83E-13
8	Capacitor	8.27E-6	1.00	8.27E-6	6.83E-13
9	PPMU	8.27E-6	1.00	8.27E-6	6.83E-13
10	Missing Solder	8.27E-6	1.00	8.27E-6	6.83E-13
11	Mfg Defects	8.27E-6	1.00	8.27E-6	6.83E-13
	Latent Failure in Phase 3	8.46E-6	1.208	1.07E-4	1.15E-10

TABLE VII
OPTIMAL CA BUDGET ALLOCATION IN PHASE 3 WITH $\theta = 0.05$

<i>i</i>	Failure Mode	c_i (\$)	b_i	x_i (\$)
1	Open Diode	430,000	1	0
2	Power Supply	150,000	1	29,996
3	EEPROM	0	0	0
4	Cold Solder	0	0	0
5	NFF	370,000	1	250,004
6	Flux Contam	0	0	0
7	SMC Limit Table	20,000	1	0
8	Capacitor	23,000	1	0
9	PPMU	310,000	1	0
10	Missing Solder	9,000	1	0
11	Mfg Defects	12,000	1	0
	Latent Failure in Phase 3	N/A	N/A	N/A

Note: if $b=c=0$, it implies CA is not considered in the decision process.

shall be applied. This is contradictory to the traditional belief that usually concentrates on the top failure modes. Given such a budget allocation scheme, the objective function would be $f(\mathbf{x}) = 4.413 \times 10^{-4}$ with $1 - \theta = 95\%$ confidence.

Table V shows the interarrival times for surfaced and latent failures observed in Phase 2. An interesting observation is that, in Phase 2, no failures occurred due to corrupted electrically erasable programmable read-only memory, cold solder, and flux contamination. Meanwhile, five latent failure modes listed from $i = 7$ to 11 occurred in Phase 2. These latent failures along with the surfaced failures will be used to predict the system failure intensity in Phase 3.

Phase 3—Continuous Improvement: In this phase, system reliability is continuously monitored and improved following

the initial CA implementation in Phase 2. Meanwhile, latent failure modes observed in Phase 2 are used to predict new latent failures in Phase 3, and the prediction result is listed in Table VI. During the prediction, it is realized that $t_c = 120\,960$ h, $T_c = 68\,544$ h, and $T = 80\,640$ h. Since the prediction is made through the end of Phase 3, hence, $t = 201\,600$ h.

New CA budget allocation can be made based on the updated failure intensities in Table VI, and the optimal budget scheme is listed in the last column of Table VII. In Phase 3, the total CA budget is $C = \$280\,000$. The new solution suggests that CA resources should be given to power supply and NFF. This is understandable because power supply is the dominant failure mode in Phase 2 and it shows an increasing trend with $\beta > 1$ by referring to Table VI. The solution also suggests that the NFF issue needs to be addressed. Although the failure quantity of open diode in Phase 2 is larger than that of NFF, no CA budget is recommended in this phase. This can be justified by the fact that the failure intensity for open diode is declining in Phase 2 with $\beta < 1$.

In this section, a fleet of ATE systems is used to demonstrate the application of the multiphase RGP model. Each phase consists of two iterative processes: reliability prediction and CA resource allocation. The optimal CA allocation in each phase is obtained using Rosen's projection algorithm. These processes are repeated until the system reliability reaches the design goal.

VII. CONCLUSION

Modeling and planning reliability growth for capital equipment is a very complex task as it involves correlated factors, such as design, manufacturing, testing, operation, and maintenance. This paper made an early attempt to seek a multiphase RGP approach taking into account the latent failures. In particular, the study proposes a CA effectiveness function and further integrates it into the reliability growth model in order to optimize the recourses against known and emerging failure modes. Each phase involves two iterative steps: reliability prediction and CA resource allocation. The optimal CA decision scheme is derived using Rosen's projection algorithm. These steps are repeated until the system reliability reaches the design target. The case study drawn from the ATE industry shows that the multiphase RGP is quite effective when products are developed in a fast time-to-market environment. As the reliability increases along with the customer shipment, more sales revenue is generated incentivizing the manufacturer to implement a broader CA program. In the future, cases drawn from other product domains should be used to verify and validate the method. The results will be compared under different design and application scenarios. The prediction of latent failures depends on the time interval specified, which may influence the CA decision making. Further analysis is anticipated in terms of identifying the best prediction period.

ACKNOWLEDGMENT

The authors would like to thank the anonymous reviewers for their constructive comments which assist the authors in improving the quality and presentation of this paper.

REFERENCES

- [1] P. Hokstad and A. T. Frøvig, "The modeling of degraded and critical failures for components with dormant failures," *Reliab. Eng. Syst. Safety*, vol. 51, no. 2, pp. 189–199, Feb. 1996.
- [2] D. S. Jackson, H. Pant, and M. Tortorella, "Improved reliability-prediction and field-reliability-data analysis for field-replaceable units," *IEEE Trans. Rel.*, vol. 51, no. 1, pp. 8–16, Mar. 2002.
- [3] J. Carlson and R. R. Murphy, "How UGVs physically fail in the field," *IEEE Trans. Robot. Autom.*, vol. 21, no. 3, pp. 423–437, Jun. 2005.
- [4] V. Grassi and S. Patella, "Reliability prediction for service-oriented computing environments," *IEEE Internet Comput.*, vol. 10, no. 3, pp. 43–49, May/Jun. 2006.
- [5] J. E. Vinson and J. J. Liou, "Electrostatic discharge in semiconductor devices: An overview," *Proc. IEEE*, vol. 86, no. 2, pp. 399–420, Feb. 1998.
- [6] T. Jin, P. Wang, and Q. Huang, "A practical MTBF estimate for PCB design considering component and non-component failures," in *Proc. Annu. Rel. Maintain. Symp.*, 2006, pp. 604–610.
- [7] X. S. Si, W. Wang, C. H. Hu, and D. H. Zhou, "Remaining useful life estimation—A review on the statistical data driven approaches," *Eur. J. Oper. Res.*, vol. 213, no. 1, pp. 1–14, Aug. 2011.
- [8] J. T. Duane, "Learning curve approach to reliability monitoring," *IEEE Trans. Aerosp.*, vol. AS-2, no. 2, pp. 563–566, Apr. 1964.
- [9] L. H. Crow, "Reliability analysis for complex, repairable systems," in *Reliability and Biometry*. Philadelphia, PA: SIAM, 1974, pp. 379–410.
- [10] M. Xie and M. Zhao, "Reliability growth plot—An underutilized tool in reliability analysis," *Microelectron. Reliab.*, vol. 36, no. 6, pp. 797–805, Jun. 1996.
- [11] C. L. Campbell, "Subsystem reliability growth allocation," in *Proc. 36th Tech. Meet., Inst. Environ. Sci.*, Mount Prospect, IL, 1990, pp. 748–751.
- [12] D. W. Coit, "Economic allocation of test times for subsystem-level reliability growth testing," *IIE Trans. Qual. Reliab. Eng.*, vol. 30, no. 12, pp. 1143–1151, Dec. 1998.
- [13] C. Benski and E. Cabau, "Unreplicated experimental design in reliability growth programs," *IEEE Trans. Rel.*, vol. 44, no. 2, pp. 199–205, Jun. 1995.
- [14] M. Krasich, J. Quigley, and L. Walls, "Modeling reliability growth in the system design process," in *Proc. Annu. Rel. Maintain. Symp.*, 2004, pp. 424–430.
- [15] M. Krasich, "Accelerated reliability growth testing and data analysis method," in *Proc. Annu. Rel. Maintain. Symp.*, 2006, pp. 385–391.
- [16] S. Inoue and S. Yamada, "Generalized discrete software reliability modeling with effect of program size," *IEEE Trans. Syst., Man, Cybern. A, Syst., Humans*, vol. 37, no. 2, pp. 170–179, Mar. 2007.
- [17] C.-G. Bai, K.-Y. Cai, Q.-P. Hu, and S.-H. Ng, "On the trend of remaining software defect estimation," *IEEE Trans. Syst., Man, Cybern. A, Syst., Humans*, vol. 38, no. 5, pp. 1129–1142, Sep. 2008.
- [18] S. Hwang and H. Pham, "Quasi-renewal time-delay fault-removal consideration in software reliability modeling," *IEEE Trans. Syst., Man, Cybern. A, Syst., Humans*, vol. 39, no. 1, pp. 200–209, Jan. 2009.
- [19] F. Xue, A. C. Sanderson, and R. J. Graves, "Multiobjective evolutionary decision support for design–supplier–manufacturing planning," *IEEE Trans. Syst., Man, Cybern. A, Syst., Humans*, vol. 39, no. 2, pp. 309–320, Mar. 2009.
- [20] T. C. Smith, "Reliability growth planning under performance based logistics," in *Proc. Annu. Rel. Maintain. Symp.*, 2004, pp. 418–423.
- [21] P. M. Ellner and J. B. Hall, "An approach to reliability growth planning based on failure mode discovery and correction using AMSAA projection methodology," in *Proc. Annu. Rel. Maintain. Symp.*, 2006, pp. 266–272.
- [22] T. Jin and H. Wang, "A multi-objective decision making on reliability growth planning for in-service systems," in *Proc. IEEE Int. Conf. Syst., Man, Cybern.*, 2009, pp. 4677–4683.
- [23] A. K. Somani, S. Palnitkar, and T. Sharma, "Reliability modeling of systems with latent failures using Markov chains," in *Proc. Annu. Rel. Maintain. Symp.*, 1993, pp. 120–125.
- [24] G. M. Susova and A. N. Petrov, "Markov model-based reliability and safety evaluation for aircraft maintenance-system optimization," in *Proc. Annu. Rel. Maintain. Symp.*, 1997, pp. 29–36.
- [25] T. Jin, H. Liao, and M. Kilari, "Reliability growth modeling for in-service systems considering latent failure modes," *Microelectron. Reliab.*, vol. 50, no. 3, pp. 324–331, Mar. 2010.
- [26] G. Pulcini, "A bounded intensity process for the reliability of repairable equipment," *J. Qual. Technol.*, vol. 33, no. 4, pp. 480–492, Oct. 2001.
- [27] L. Attardi and G. Pulcini, "A new model for repairable systems with bounded failure intensity," *IEEE Trans. Rel.*, vol. 54, no. 4, pp. 572–582, Dec. 2005.
- [28] P. Wang and D. W. Coit, "Repairable systems reliability trend tests and evaluation," in *Proc. Annu. Rel. Maintain. Symp.*, 2005, pp. 416–421.

- [29] T. Jin, Y. Yu, and F. Belkhouche, "Reliability growth using retrofit or engineering change order—A budget-based decision making," in *Proc. Ind. Eng. Res. Conf.*, 2009, pp. 2152–2157.
- [30] S. Ross, *A First Course in Probability*, 8th ed. Englewood Cliffs, NJ: Prentice-Hall, 2009.
- [31] J. B. Rosen, "The gradient projection method for non-linear programming, Part I," *J. Soc. Ind. Appl. Math.*, vol. 8, no. 1, pp. 181–217, Mar. 1960.
- [32] M. S. Bazaraa and C. M. Shetty, *Nonlinear Programming: Theories and Applications*, 3rd ed. New York: Wiley, 2006.
- [33] Y.-S. Dai, M. Xie, and X. Wang, "A heuristic algorithm for reliability modeling and analysis of grid systems," *IEEE Trans. Syst., Man, Cybern. A, Syst., Humans*, vol. 37, no. 2, pp. 189–200, Mar. 2007.
- [34] A. Sutcliffe, W.-C. Chang, and R. S. Neville, "Applying evolutionary computing to complex systems design," *IEEE Trans. Syst., Man, Cybern. A, Syst., Humans*, vol. 37, no. 5, pp. 770–779, Sep. 2007.
- [35] L. H. Crow, "Confidence interval procedures for the Weibull process with applications to reliability growth," *Technometrics*, vol. 24, no. 1, pp. 67–72, Feb. 1982.



Tongdan Jin (S'97–M'01) received the M.S. degree in electrical and computer engineering and the Ph.D. degree in industrial and systems engineering from Rutgers University, Camden, NJ.

He is an Assistant Professor with the Ingram School of Engineering, Texas State University, San Marcos. His research has been published in Reliability Engineering and Systems Safety and Microelectronics Reliability, among others. His research interests include system reliability modeling and optimization, electronics prognostics and diagnostics,

and performance-based logistics management.



Ying Yu received the Ph.D. degree in systems engineering from Southeast University, Nanjing, China, in 2009.

Between November 2008 and May 2009, she was a Visiting Scholar with Texas A&M International University, Laredo. She is with the Department of Automation, Shanghai University, Shanghai, China. She has published papers in peer-reviewed journals and frequently presented works in academic conferences. Her research interests include fuzzy/probabilistic supply chain modeling,

simulation, and optimization.



Hong-Zhong Huang received the Ph.D. degree in reliability engineering from Shanghai Jiao Tong University, Shanghai, China.

He is a Full Professor and the Dean of the School of Mechatronics Engineering, University of Electronic Science and Technology of China, Chengdu, China. He has published 150 journal papers and five books in the field of reliability engineering. He has held visiting appointments at several universities in the USA, Canada, and Asia. His current research interests include system reliability analysis, warranty,

maintenance planning and optimization, and computational intelligence in product design.

Dr. Huang was the recipient of the Golomski Award from the Institute of Industrial Engineers in 2006.