

Designing Guardbands for Instantaneous Aging Effects

Victor M. van Santen¹, Hussam Amrouch¹, Javier Martin-Martinez²,
Montserrat Nafria² and Jörg Henkel¹

¹ Karlsruhe Institute of Technology, ² Universitat Autònoma de Barcelona
{victor.santen, amrouch, henkel}@kit.edu
{javier.martin.martinez, montse.nafria}@uab.cat

ABSTRACT

Bias Temperature Instability (BTI) is one of the key causes of reliability degradations of nano-CMOS circuits. While the long-term impact of BTI has been studied since years, the *short-term* implications of BTI on circuits are unexplored. In fact, in physics short-term BTI effects, i.e. instantaneous (i.e. sub μs) frequency dependent processes, have been recently reported. In order to design circuits with guardbands that are safe for long-term *and* instantaneous effects, new aging models are required. We are presenting the first approach that in fact considers both long-term as well as instantaneous BTI effects. It can be employed for complex circuits at the micro-architecture level. Designing guardbands based upon our physical BTI model reduces the guardbands by 41% and thus allows for the development of more cost-effective yet reliable designs. We also revisit existing state-of-the-art aging mitigation techniques to investigate how they can be properly adapted to additionally account for instantaneous aging effects. Along with our BTI model this further reduces the guardbands by up to 59%.

Download Aging Estimation: This work is publicly available at <http://ces.itec.kit.edu/dependable-hardware.php>

1. INTRODUCTION

Modeling and mitigating aging effects are key challenges of this decade since reliability must not be compromised, while the current nano-CMOS is highly susceptible to aging. Bias Temperature Instability (BTI) is recognized as one of the major aging phenomena due to its considerable ability to degrade the electrical characteristics of MOSFETs. To sustain reliability, aging degradations need to be estimated at design time in order to provide the required guardband (i.e. designing the system above specification to tolerate degradation) that protects circuits against aging effects.

The major challenge is that BTI-induced degradations are estimated solely regarding its well-known *long-term* impact. The implications of *short-term* BTI on circuits are unexplored. In fact, reliability physics report that BTI consists of instantaneous (i.e. sub μs) frequency dependent processes, which were uncovered due to advances in measurement tech-

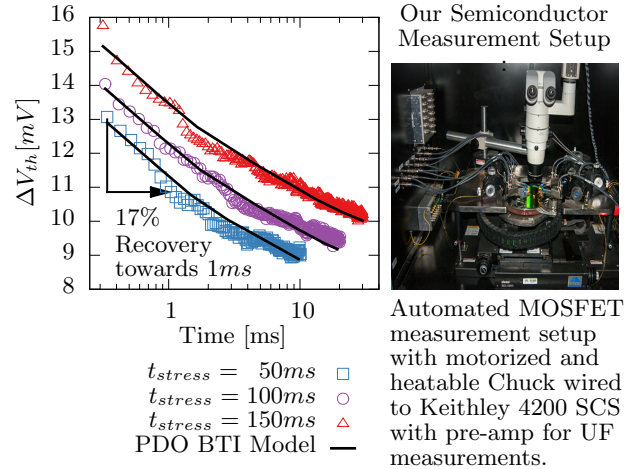


Figure 1: BTI recovery, measured on our ultra-fast measurement equipment, validating the short-term behavior of the base BTI model [1] we rely upon in our implementation. Additionally, our measurement highlights the importance of UF measurements as BTI recovers 17% degradation from $0.34ms \rightarrow 1ms$ after stress. Indicating how slow measurement missed the actual impact of BTI.

niques [2] (explained in detail in section 2). The impact of instantaneous BTI is considerable and such instantaneous shifts may suddenly violate the employed guardbands manifesting itself as BTI-induced errors.

In order to prevent guardband violations due to instantaneous BTI, a new aging estimation approach is required that considers instantaneous and long-term BTI jointly to design guardbands protecting against both. There are two major challenges: First, guardbands cannot be further increased to incorporate new BTI effects. Actually, with each new technology generation the available design space for guardbands shrinks as different sources of reliability degradation phenomena (aging, noise, process variation, etc.) increase, while the resiliency against them decreases with the decrease in supply voltage [3]. Containing guardbands within available design space (i.e. designing *narrow guardbands*), requires accurate models along with a design methodology that replaces worst-case assumptions with actual occurring aging to minimize overestimation. The second challenge is that BTI-induced degradations must be estimated for instantaneous and long-term BTI at the micro-architecture level, which requires a fast, i.e. computational lightweight model to model such complex circuitry in feasible simulation times.

Diverse approaches for BTI modeling exist ranging from the physical level [1] towards the micro-architecture level [4]. At the physical level BTI is measured based upon defect concentrations in transistors and its impact is expressed

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. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

DAC '16, June 05-09, 2016, Austin, TX, USA

© 2016 ACM. ISBN 978-1-4503-4236-0/16/06...\$15.00

DOI: <http://dx.doi.org/10.1145/2897937.2898006>

as induced shifts in transistor parameters (threshold voltage shift ΔV_{th}) [1]. These BTI models, which model the underlying physical processes of BTI to estimate it, are *physical BTI models*. In contrast, at the micro-architecture level, BTI is measured by observing failure rates of chips over time. Then BTI is expressed by simple equations fitted to mimic the observed failure rates in simulations which model shifts in transistor parameters (ΔV_{th}) [4]. BTI models with equations fitted to match chip failure behavior are called *empirical BTI models* in this work.

Interestingly, the *physical* and *empirical* approach differs significantly due to the direct (transistor degradation) and indirect (chip failure rates) calibration with measurements. Empirical models have a high degree of uncertainty due to the probabilistic nature of chip failures [4]. To ensure reliable designs, circuit designers must consider the worst samples of these distributions and design their guardbands accordingly.

Despite their inherent uncertainty, *empirical BTI models* are used as their simplicity and speed allows BTI estimations within complex circuitry. However, to carefully design *narrow guardbands*, the *physical models* are more suitable as their detailed modeling reduces uncertainty, providing results closer to the actually required guardbands. Therefore, commercial design tools like MOSRA from Synopsys employed *physical models* [5]. Since physical models are computational infeasible, MOSRA reduced the number of mathematical terms in their hot carrier model to limit computational and calibration complexity at the cost of compromising accuracy [5]. Despite those efforts simplifying *physical models*, MOSRA is only applicable to circuits with moderate complexity [6]. Academia also attempts to solve the performance problem, [7] reduced data which needs to be processed and [6] employs an offline look up table approach. Unfortunately, neither approach is sufficiently fast as [7] still calculates thousands of data points for a single transistor, while the look up tables for [6] can become unfeasible for complex circuitry evaluated over a wide range of operating conditions. Left without feasible *physical BTI modeling*, circuit designers are forced to employ *empirical models* despite their overestimation.

Our novel contributions:

1. We present a *physical BTI model* incorporating both instantaneous *and* long-term effects of BTI. It is computationally lightweight to be feasible estimating complex circuitry, while it retains the accuracy of physical models.
2. Adapting existing aging mitigation techniques to reduce the guardbands further by reducing the stimuli of long-term and instantaneous aging.

2. INSTANTANEOUS BTI

BTI is stimulated by transistor activity, i.e. the transistor degrades when it is on (i.e. under stress for time t_{stress}) and recovers when it is in off-state (i.e. in recovery for $t_{recovery}$). Activity waveforms, i.e. series of t_{stress} & $t_{recovery}$, can be summarized with the on-/off-ratio λ and the frequency of the state changes.

2.1 Exposing Instantaneous BTI

Traditionally, BTI is measured with the measure-stress-measure (MSM) pattern, i.e. stressing the device for t_{stress} and then removing the stress for $t_{measure}$ (from 1ms to 1s) from the device to characterize the device parameters (e.g. ΔV_{th}) [2]. Then ultra-fast (UF) measurement techniques were introduced [8] in which $t_{measure} < 1ms$. UF

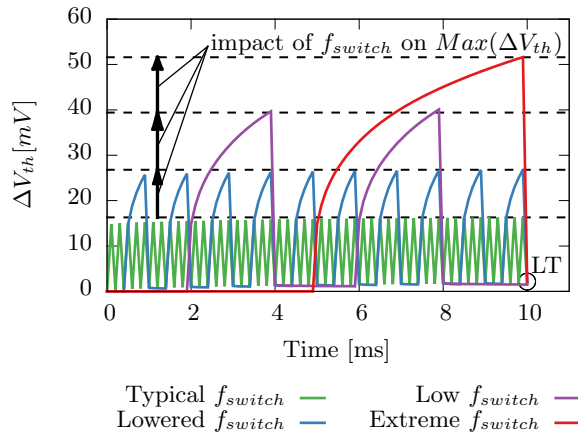


Figure 2: Frequency Dependence of BTI ($\lambda = 0.5, T = 80^\circ C, V_{dd} = 1.0V$) with rising $Max(\Delta V_{th})$ for lower f_{switch} . Frequency independent long-term point marked with "LT". f_{switch} based upon Fig 5.

measurements uncovered that BTI does not solely accumulate over time, degrading reliability, but additionally reacts *instantaneously* (i.e. sub-microsecond) to stimuli with degradation or recovery. In fact, after the stress was removed in MSM, BTI partially recovered from high levels of degradations altering the perception of BTI from its actual instantaneous nature. Fig 1 shows recovery below 1ms measured by UF measurements highlighting how 17% degradation may be missed when $t_{measure} = 1ms$.

2.2 Frequency Dependency of Instantaneous BTI

We differentiate two frequencies, the operating frequency $f_{operation}$ of the circuit (i.e. the clock frequency) and the switching frequency f_{switch} of individual transistors (i.e. the frequency of transistor switches from on to off states).¹ Note, f_{switch} is only loosely coupled with $f_{operation}$. For example, a memory cell storing the same data or clock gated logic, does not switch ($f_{switch} = 0Hz$) regardless of $f_{operation}$. Therefore, most significant bits of memory or an ALU, storing same values for prolonged times [7], switch infrequently ($f_{switch} \ll f_{operation}$).

Long-term BTI is frequency independent, i.e. multiple waveforms with identical λ but varying f_{switch} lead to identical results [7]. With identical λ , the stress/recovery ratio is fixed and e.g. longer stress phases are compensated by longer recovery phases, ultimately resulting in the same BTI-induced degradation.

Instantaneous BTI, on the other hand, is frequency dependent [9], [10]. As the physical processes of BTI are in reality instantaneous, individual stress phases itself lead to considerable degradations. Instantaneous recovery cannot compensate, if a single instantaneous stress phase already violates the guardband.

Fig. 2 shows the BTI-induced degradation ΔV_{th} for three different f_{switch} but identical λ . Instantaneous BTI exhibits degradation peaks after each stress phase. Intolerable degradation ($Max(\Delta V_{th})$) is reached at low f_{switch} , forcing the consideration of instantaneous BTI. Other observations are, that degradation decreases at higher frequencies (also see Fig. 3) like [10], [9] reported. At the same time, long-term BTI degradation (marked with "LT") exhibits the same ΔV_{th}

¹Note, that f_{switch} is different from the number of switches (toggling rate) used for hot carrier modeling.

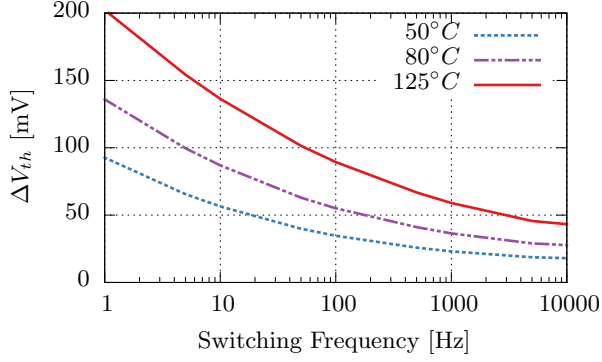


Figure 3: ΔV_{th} due 1000s with $\lambda = 0.5$ and 1 period continuous stress with $t_{stress} = \frac{1}{f_{switch}}$ for 3 temperatures.

when all frequencies are in phase, i.e. being f_{switch} independent as [7] claimed.

Note, that claims like frequency independence due to high $f_{operation}$ [5] are incorrect, as f_{switch} in memory cells, clock gated logic or the most significant bits switch at significantly lower frequencies than $f_{operation}$, i.e. $f_{switch} \ll f_{operation}$ (see Fig. 5).

3. OUR PROPOSED BTI MODEL

In order to estimate guardbands for circuit designs, a BTI model must consider long-term *and* instantaneous BTI, be computationally lightweight for the feasible employment in complex circuitry, while having a low uncertainty for *narrow guardbands*. We therefore, enhanced the physical model from [1] to directly calculate maximum degradation $Max(\Delta V_{th})$ based upon λ and f_{switch} instead of waveforms with stress time t_s / recovery time t_r . Designing guardbands requires solely $Max(\Delta V_{th})$ occurring during the desired lifetime of the circuit (t_{life}), so the tedious calculation of ΔV_{th} over time could be removed. The model is able to predict our UF measurements well (see Fig. 1) which validated that it can model instantaneous BTI. Both the original and reshaped model calculate ΔV_{th} by integrating over the defect distribution D and the occupancy probability P_{occ} of the defects. The latter was replaced in our implementation, to calculate only $Max(\Delta V_{th})$.

3.1 Original BTI Model

The probabilistic defect occupancy model (PDO) [1] calculates ΔV_{th} for a given temperature T , a voltage V , stress time t_s and recovery time t_r :

$$\Delta V_{th}(t_s, t_r) = N \cdot \bar{\eta} \int_0^\infty \int_0^\infty D(\tau_e, \tau_c) \cdot P_{occ}(\tau_c, \tau_e; t) d\tau_e d\tau_c \quad (1)$$

$$\text{with } \tau_c = \tau_c(T, V) \text{ and } \tau_e = \tau_e(T, V)$$

BTI is modeled by integrating over two distributions. First $D(\tau_e, \tau_c)$ as the defect distribution, i.e. the distribution of defects characterized with their carrier capture τ_c and emission τ_e times. This characterization of the defect distribution is performed with measurements of the gate dielectric at different T, V as τ_c and τ_e are dependent on the temperature T and voltage V applied to the transistor [7].

The second distribution is the occupancy map P_{occ} , i.e. the occupation probability of a defect given by the current and past activity of the transistor. For a given stress time t_s , all defects with $\tau_c < t_s$ have likely captured a carrier and therefore contribute with their ΔV_{th} towards the overall

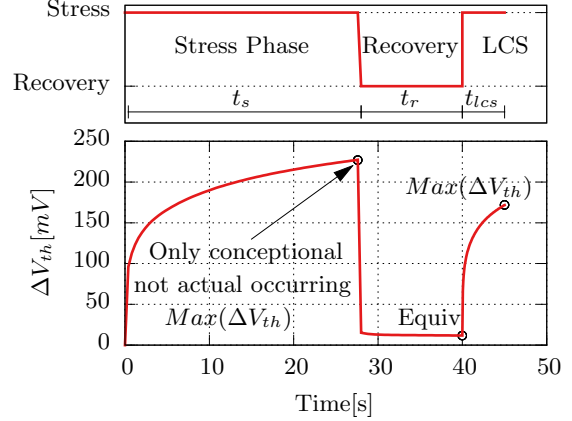


Figure 4: Top Plot: Stress and recovery states are illustrated annotated with stress phase t_s , recovery phase t_r and longest continuous stress (LCS) phase t_{LCS} . Bottom plot: Corresponding ΔV_{th} with the equivalent long-term point marked "Equiv" and the final result marked " $Max(\Delta V_{th})$ ". Note: Unrealistic values used in plot to show principle, as in actual systems $t_{LCS} \ll (t_s, t_f)$. For actual results see Fig. 6.

ΔV_{th} . Then for a given recovery time t_r , all defects currently occupied due previous stress phases and $\tau_e < t_r$ likely release their carrier, i.e. do not contribute to the overall ΔV_{th} any more. According to [1] P_{occ} for digital voltage waveforms can be expressed as:

Stress:

$$P_{occ}(t) = P_{occ}(t_i) + \left(\frac{\tau_e}{\tau_e + \tau_c} - P_{occ}(t_i) \right) \cdot \left(1 - e^{-\frac{t_i - t}{\tau_{sr}}} \right) \quad (2)$$

Recovery:

$$P_{occ}(t) = \frac{\tau_e}{\tau_e + \tau_c} + \left(P_{occ}(t_i) - \frac{\tau_e}{\tau_e + \tau_c} \right) \cdot \left(e^{-\frac{t_i - t}{\tau_{sr}}} \right) \quad (3)$$

$$\text{with } \tau_{sr} = \frac{1}{\frac{1}{\tau_e} + \frac{1}{\tau_c}} \quad \tau_c = \tau_c(T, V) \quad \tau_e = \tau_e(T, V)$$

with t_i the time of the i -th switch between on and off state.

Integrating over both distributions from 0 to ∞ , captures the impact of all occupied defects (defects capturing a carrier), i.e. contributing their ΔV_{th} to the overall ΔV_{th} .

3.2 Long-Term Phases

In our estimation of BTI-induced $Max(\Delta V_{th})$, we split the calculation in two parts. First we estimate the long-term BTI-induced degradation for the desired lifetime t_{life} of the circuit. Then in a second step, we consider the degradation due to instantaneous BTI on top of long-term BTI.

Section 2.2 explained how long-term BTI is frequency independent if λ remains identical. Therefore, we can model the history of the transistor with overall λ for the entire lifetime and then we can replace the waveform with solely two data points: First a *stress phase* for time t_s , then a *recovery phase* for t_r :

$$t_s = \lambda \cdot t_{life} \quad (4)$$

$$t_r = (1 - \lambda) \cdot t_{life} \quad (5)$$

3.3 Longest Continuous Stress Phase

To account for the instantaneous effects of BTI, we introduce the *longest continuous stress (LCS) phase*. After the *stress* and *recovery phase*, we stress the transistor for

longest continuous stress occurring in the activity waveform ($Max(t_{stress})$):

$$t_{lcs} = Max(t_{stress}) \quad (6)$$

$$\text{or } t_{lcs} = t_{period}(Min(f_{switch})) = \frac{1}{Min(f_{switch})} \quad (7)$$

Using the longest continuous stress, ensures we catch the worst shift due to instantaneous BTI as longer stress result in higher ΔV_{th} (see Fig. 3). By placing the frequency phase at the end, we mimic the occurrence of the longest continuous stress phase at t_{life} , i.e. the worst instantaneous shift occurs on top of the worst long-term degradation ensuring that guardbands can tolerate long-term and instantaneous BTI jointly. As most computing systems are periodical, i.e. execute the same tasks regularly, $Max(t_{stress})$ occurs periodically. Therefore $Max(t_{stress})$ occurs also at end of lifetime t_{life} , indicating that our worst-case assumption is not an overestimation in most computing systems.

The final activity waveform is shown in Fig. 4. On the top, an activity waveform for $\lambda = 0.7$, $t_{life} = 40s$, $f_{switch} = 0.2Hz \rightarrow Max(t_{stress}) = 5s$, $T = 80^\circ C$, $V = 1.0V$ is shown, resulting in $t_s = 28s$, $t_r = 12s$ and $t_f = 5s$.

Note that, the end-point in the bottom plot $Max(\Delta V_{th})$ is not the highest plotted ΔV_{th} . Stress phase t_s is solely a concept to reduce computational complexity and not an actual occurring stress phase.

3.4 Simplified Occupancy Probability

Instead of calculating P_{occ} for an arbitrary waveform, P_{occ} is now calculated for 2 defined transitions between the 3 phases, i.e. stress \rightarrow recovery \rightarrow stress at known time-steps ($t_s \rightarrow t_r \rightarrow t_{lcs}$). These defined inputs allow us to simplify the P_{occ} to a 3 step calculation:

1. Stress Phase:

$$P_{occ}(t_s) = \left(\frac{\tau_e}{\tau_e + \tau_c} \right) \cdot \left(1 - e^{-\frac{t_s}{\tau_{sr}}} \right) \quad (8)$$

2. Recovery Phase:

$$P_{occ}(t_r) = \frac{\tau_e}{\tau_e + \tau_c} + \left(P_{occ}(t_s) - \frac{\tau_e}{\tau_e + \tau_c} \right) \cdot \left(e^{-\frac{t_s - t_r}{\tau_{sr}}} \right) \quad (9)$$

3. Longest Continuous Stress Phase:

$$P_{occ}(t_{lcs}) = P_{occ}(t_r) + \left(\frac{\tau_e}{\tau_e + \tau_c} - P_{occ}(t_r) \right) \cdot \left(1 - e^{-\frac{t_r - t_{lcs}}{\tau_{sr}}} \right) \quad (10)$$

$$\text{with } \tau_{sr} = \frac{1}{\frac{1}{\tau_e} + \frac{1}{\tau_c}} \quad \tau_c = \tau_c(V, T) \quad \tau_e = \tau_e(V, T)$$

Originally P_{occ} is calculated recursively for every state switch in the waveform and a recursion depth of #switches (e.g. 1,467,315 for an average transistor while the processor executes ‘‘barnes’’). In contrast, our new P_{occ} calculation takes exactly 3 steps with a recursion depth of 3:

$$P_{occ} = P_{occ}(t_{lcs}, P_{occ}(t_r, P_{occ}(t_s))) \quad (11)$$

Simplifying P_{occ} results in a significant speed-up as P_{occ} is updated every time the voltage or temperature changes to account for the temperature and voltage dependence of $\tau_e(V, T)$, $\tau_c(V, T)$.

As ΔV_{th} is a function of P_{occ} (see eq. 1) we obtain:

$$\Delta V_{th}(\lambda, f_{switch}, T, V, t_{life}) = N \cdot \bar{\eta} \int_0^\infty \int_0^\infty D \cdot P_{occ} d\tau_e d\tau_c \quad (12)$$

The model provides the maximum BTI-induced degradation $Max(\Delta V_{th})$ for given operating conditions (λ , f_{switch} , T , V , t_{life}) in milliseconds, as just 3 phases are processed, while still employing detailed modeling of physical processes for the calculation of instantaneous and long-term BTI.

4. GUARDBAND ESTIMATION

Estimating the guardbands for a circuit requires input parameters for the BTI model to estimate the degradation. Worst-case scenarios are the safe and easy option, i.e. assuming the highest T , worst λ , slowest f_{switch} , etc. However, this leads to guardbands which exceed the available design space.

The alternative is to estimate the aging stimuli based upon workload of the circuit, modeling actual occurring aging [7]. To estimate the activity of the workload, we employ gem5 [11] as cycle accurate system simulator, which simulates the execution of the workload in Linux 2.6 on ALPHA 21264 Out-of-Order processor at $f_{operation} = 2GHz$.

In the following section, we exemplify our approach on the register file of our microprocessor. The implementation is not limited to register files and monitors other microprocessor components (caches, ALUs, etc.) in a analogous manner.

Guardband: In our scenario, the guardband is defined as the BTI-induced degradation in percent of the static noise margin of the SRAM cells in a register [12],[13].

Activity Monitoring: We implemented our own architecture level activity monitor, which estimates the signal probabilities (ratio for 0 or 1) for each bit within the register file. Assuming an SRAM-based register file, we determine the activity based upon the signal probabilities of the storage bits and addressing bits to calculate λ for each transistor within the cells.

Frequency Monitoring: To obtain f_{switch} our activity monitor monitors the longest, shortest and average time for the data stored in the SRAM cells. The longest period defines $Min(f_{switch})$ which is of main interest.

Power Estimation: To estimate the power of the processor, we employ McPat, [14], which models the static and dynamic power consumption of the ALPHA processor. McPat takes the gem5 activity waveform of each microprocessor component and translate it to power waveforms for each component.

Temperature Estimation: The power waveform for each microprocessor component is passed to the thermal simulator HotSpot. [15]. Together with the floorplan of the microprocessor its temperature can be estimated. Power and temperature estimation is iterative, as the temperature recursively depends on the static power consumption (leakage) of the processor.

These time consuming steps are performed once, as these aging stimuli do not change as long as the micro-architecture of the processor remains identical.

4.1 Designing Guardbands

Our workload monitoring provides λ , $Min(f_{switch})$ and T for each individual workload. Together with t_{life} , V given by the specification, we can estimate $Max(\Delta V_{th})$ for the given workload with our proposed BTI model.

With $Max(\Delta V_{th})$ known, degraded SPICE simulations of the register file can be performed to estimate if the register file operates within specification. In practice, this means employing Monte Carlo SPICE simulations of SRAM cells to model variability introduced by manufacturing together with aging-induced degradation to verify meeting time constraints (e.g. read access time < clock period) and sufficient resiliency against noise (static noise margin) or radiation (critical charge). A circuit designer can then adapt the

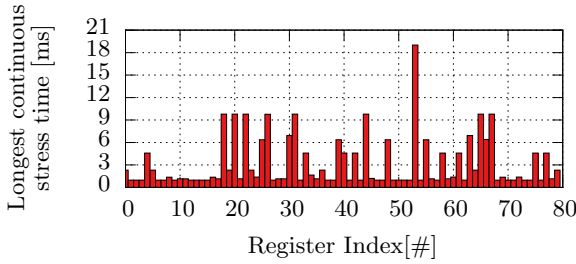


Figure 5: Longest time the same value is stored in a register of the register file while executing the “barnes” application at $f_{operation} = 2GHz$. These LCS phases stimulate instantaneous BTI and are the input for the calculation of $Max(\Delta V_{th})$.

guardbands until probability of failure of the circuit P_{fail} is below P_{fail} of the specification. Once the smallest guardband is found, which still satisfy the specification, the circuit designer must employ the guardband in his circuit. For example, he could up-size the transistor widths to make transistors faster and more resilient against noise or reduce the desired operating frequency to increase the timing slack.

5. EVALUATION

Model Validation: Simplifying the model did introduce only minor inaccuracies, due long-term BTI not being perfectly frequency independent and $Max(t_{stress})$ not occurring exactly at t_{life} . Compared to the carefully validated PDO model [1], which results in $Max(\Delta V_{th}) = 73.12mV$ for low f_{switch} , $T = 125^{\circ}C$, $V = 1.5V$, $\lambda = 0.5$, $t_{life} = 10$ years, our model estimates $Max(\Delta V_{th}) = 72.89mV$, i.e. a deviation of 0.3%. The $Max(\Delta V_{th})$ deviation between PDO and our model for the experiments in Fig. 6 was below 0.3%.

Model Performance: To evaluate the performance of the model, we generated an activity waveform with 10^6 data points and compared original PDO [1], PDO with compressed (10^4) waveform [7] and the proposed model in this work. Our model required 0.094s to estimate $Max(\Delta V_{th})$, while PDO needed 4.366s with and 433.5s without compression, i.e. speedup of 99x compared against [7] and 4567x against [1]. Note, that the execution times of our model and the compression are almost independent of the waveform length (waveform analysis depends on #points, then computation on fixed #points), while PDO has an execution time proportional to the waveform size. Our model could perform the aging estimation for a DCT circuit featuring 350,015 transistors in 9.2 hours and IDCT circuit in 9.1 hours. DCT-IDCT circuits are often employed in image processing and are $\sim 3x$ larger than a typical RISC processor[16], highlighting how our physical model is feasible for complex circuitry at micro-architecture level. For designs exceeding this level of complexity, our approach in [16] can be used jointly with the BTI model presented in this work.

Intra-Application-Variation: Fig. 5 highlights the intra-application variation of the instantaneous BTI stimulant $Max(t_{stress})$ across the register file. Most registers change their state regularly during the execution of the “barnes” application, i.e. hold their state not longer than $2ms \rightarrow f_{switching} = 500Hz$. Register 53 exhibits large $Max(t_{stress})$ as the same state is held for almost 20ms which results in $f_{switching} = 50Hz$. Slow switching registers like 53 exhibit $Max(\Delta V_{th}) \approx 100mV$, while the average registers exhibit 63mV and the best just 54mV.

Inter-Application-Variation: While a single application shows different behavior among its registers, the applications themselves are similar to each other in terms of state changes, i.e. the worst and average $Max(t_{stress})$ in the entire register file is almost identical for each application. We

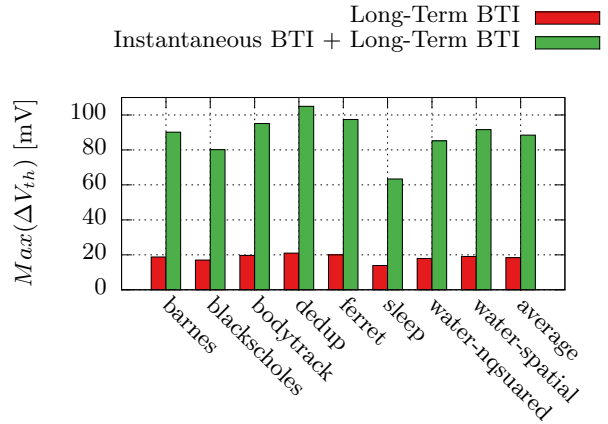


Figure 6: $Max(\Delta V_{th})$ for applications at their corresponding operating conditions ($V = 1.0V$, $T \in [49^{\circ}C, 85^{\circ}C]$). Even though $Max(t_{stress})$ is similar for all applications different λ, T lead to different $Max(\Delta V_{th})$.

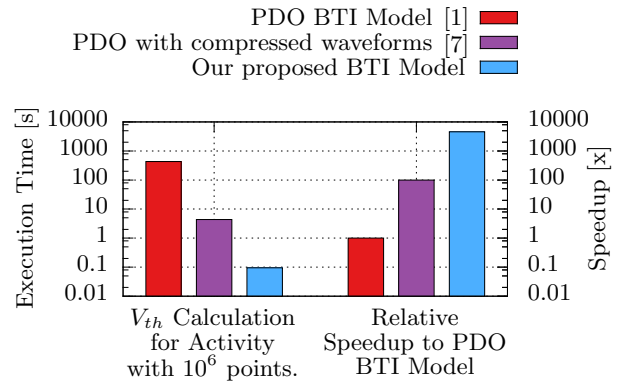


Figure 7: Comparison of BTI models, which base themselves on the PDO BTI Model [1]. Our proposed model has 94ms execution time, resulting in a speedup of 4567x compared to PDO.

assume that the register renaming of the out-of-order alpha processor is the main reason for this observation. Even though $Max(t_{stress})$ may be similar for the applications, other operating conditions like T, λ are not. Hence the induced $Max(\Delta V_{th})$ shown in 6 is different, highlighting that the joint impact between the operating conditions must be considered.

Impact of Instantaneous BTI: In Fig. 6 the impact of long-term BTI versus long-term BTI with instantaneous BTI is shown, illustrating that instantaneous BTI contributes 70.1mV or 79,2% to the overall average $\Delta V_{th} = 88.4mV$. This motivates mitigating instantaneous BTI as it predominantly governs BTI guardbands.

Impact on Guardbands: In Fig. 8 the guardbands of our register file scenario are illustrated for aging stimuli based upon the average of the studied benchmarks. Designing the guardbands with an empirical model leads to 8%, while our approach estimates 4.66% to be sufficient to tolerate the occurring aging-induced degradation. The guardband reduction due to mitigation via periodic inversion is discussed in the next section.

6. ADAPTING EXISTING AGING MITIGATION TECHNIQUES

As Fig. 6 shows, the additional guardband required to tolerate instantaneous BTI is considerable. In order to design narrow guardbands mitigation techniques are necessary, to reduce the required guardband. These mitigation techniques

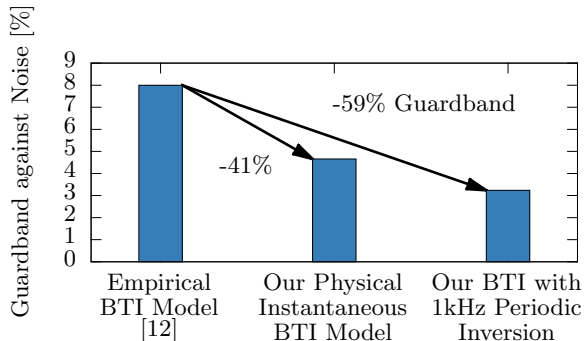


Figure 8: Reduction of the required guardband to protect against data corruption due to noise in our register file scenario, due to the estimation of the actual required guardband. Further reduction due to employed aging mitigation technique.

should optimize long-term and instantaneous BTI jointly, as both degradations define the required guardband.

Thermal Management/Voltage Scaling: Techniques focusing on reducing a single aging stimulant like thermal management (reduce T) or voltage scaling (reduce V), must take instantaneous aging effects into account. Their policies were designed without f_{switch} in mind as they were evaluated with *empirical models* (i.e. neglecting frequency dependence of BTI) potentially resulting in strong stimulation of instantaneous BTI. For example, thermal management can stall activity (clock gating) to reduce the dynamic power consumption of the circuit, which directly affects f_{switch} . Similarly, when voltage scaling reduces V to save energy or limit generated heat, $f_{operation}$ is lowered to prevent timing errors, which in turn decreases f_{switch} if the same workload is executed. These side-effects must be considered when employing such mitigation techniques, i.e. updating their policies using our proposed BTI model to find the pareto-optimal solution considering all operating conditions jointly including f_{switch} for instantaneous BTI.

Periodic inversion: Originally intended to distribute aging stress as uniformly as possible within circuits, this technique inverts logic signals periodically [17]. Circuits are extended with a flag, indicating if the data is currently *normal* or *inverted*. In *normal* mode everything is regularly processed, while in *inverted* mode data is either processed in an inverted manner (knowing that the result is also inverted) or temporarily returned to its original state during processing. Periodically inverting the entire system ensures that transistors operate close to $\lambda = 0.5$ [17], which reduces aging stimuli for transistors which had $\lambda > 0.5$ while it increases stimuli for $\lambda < 0.5$. This reduces the variability, raising the lower boundary for λ and hence reducing the guardband.

To have a first order approximation for the overhead, we refer to our register file scenario in the Alpha 21264. The 80 registers in the register file would require 1 read and 1 write operation every $1ms$ to invert the data (read, invert, write back) with $f_{inversion} = 1kHz$. With ≈ 25000 reads and ≈ 10000 writes per $1ms$ at $f_{operation} = 2GHz$ across our studied benchmarks the performance and power overhead would be negligible. This rough estimation supports [17] claiming less than 1% performance impact if logic and caches are protected at $t_{period}(inversion) = 1ms$.

Next to its original intent, periodic inversion can be employed to ensure f_{switch} does not fall below a lower boundary. When the logic signals are inverted, the states of all transistors change ensuring $f_{switch} \geq f_{inversion}$. Inverting every $1ms$, i.e. $f_{switch} \geq 1kHz$ would reduce the average $Max(\Delta V_{th})$ across the studied benchmarks by $52mV$, very close the long-term only result. The employment of periodic inversion together with our BTI model lead to a reduction of 59% of the guardband in our register file scenario.

7. CONCLUSION

We presented the first physical BTI model that models the frequency dependent instantaneous effect of BTI at the micro-architecture level. Providing operating waveforms as aging stimuli to the physical BTI model enables designing narrow guardbands. Additionally, existing aging mitigation techniques are adapted to reduce the stimuli for long-term and instantaneous BTI, resulting in up to 59% guardband reduction.

Acknowledgements

We would like to acknowledge Michael Skinder for his assistance in performing the experiments. This work was supported in parts by the German Research Foundation (DFG) as part of the priority program “Dependable Embedded Systems” [18] (SPP 1500 - spp1500.itec.kit.edu).

The UAB group acknowledges the support of the Spanish MINECO, ERDF (TEC2013-45638-C3-1-R) along with the Generalitat de Catalunya (2014SGR-384).

8. REFERENCES

- [1] J. Martin-Martinez, B. Kaczer, M. Toledano-Luque *et al.*, “Probabilistic defect occupancy model for NBTI,” in *IRPS*, 2011.
- [2] H. Reisinger, O. Blank, W. Heinrigs *et al.*, “Analysis of nbtI degradation- and recovery-behavior based on ultra fast vt-measurements,” in *IRPS*, 2006.
- [3] K. Aadithya, A. Demir, and S. Venugopalan *et al.*, “Accurate Prediction of Random Telegraph Noise Effects in SRAMs and DRAMs,” *TCAD*, 2013.
- [4] J. Henkel, T. Ebi, H. Amrouch *et al.*, “Thermal management for dependable on-chip systems,” in *ASP-DAC*, 2013.
- [5] B. Tudor, J. Wang, C. Sun *et al.*, “MOSRA: An efficient and versatile MOS aging modeling and reliability analysis solution for 45nm and below,” in *ICSICT*, 2010.
- [6] J. Chen, S. Wang, and M. Tehranipoor, “Critical-reliability Path Identification and Delay Analysis,” *J. Emerg. Technol. Comput. Syst.*, 2014.
- [7] H. Amrouch, J. Martin-Martinez, V. van Santen *et al.*, “Connecting the physical and application level towards grasping aging effects,” in *IRPS*, 2015.
- [8] H. Reisinger, U. Brunner, W. Heinrigs *et al.*, “A Comparison of Fast Methods for Measuring NBTI Degradation,” *TDMR*, 2007.
- [9] T. Grasser, B. Kaczer, H. Reisinger *et al.*, “On the frequency dependence of the bias temperature instability,” in *IRPS*, 2012.
- [10] S. Mahapatra, N. Goel, S. Desai *et al.*, “A Comparative Study of Different Physics-Based NBTI Models,” *T-ED*, 2013.
- [11] N. Binkert, B. Beckmann, G. Black *et al.*, “The Gem5 Simulator,” *SIGARCH Comput. Archit. News*, 2011.
- [12] S. Kumar, C. Kim, and S. Sapatnekar, “Impact of nbtI on sram read stability and design for reliability,” in *ISQED*, 2006.
- [13] H. Amrouch, V. van Santen, T. Ebi *et al.*, “Towards interdependencies of aging mechanisms,” in *ICCAD*, 2014.
- [14] S. Li, J. H. Ahn, R. D. Strong *et al.*, “The McPAT Framework for Multicore and Manycore Architectures: Simultaneously Modeling Power, Area, and Timing,” *ACM Trans. Archit. Code Optim.*, 2013.
- [15] M. R. Stan, K. Skadron, M. Barcella *et al.*, “Hotspot: a dynamic compact thermal model at the processorarchitecture level,” *Microelectronics Journal*, 2003.
- [16] H. Amrouch, B. Khaleghi, A. Gerstlauer *et al.*, “Reliability-Aware Design to Suppress Aging,” in *DAC*, 2016.
- [17] E. Gunadi, A. A. Sinkar, N. S. Kim *et al.*, “Combating Aging with the Colt Duty Cycle Equalizer,” in *MICRO*, 2010.
- [18] J. Henkel, L. Bauer, J. Becker *et al.*, “Design and architectures for dependable embedded systems,” in *CODES*, 2011.